



PLS: a new statistical insight through the prism of orthogonal polynomials

Mélanie Blazère, Fabrice Gamboa, Jean-Michel Loubes

► To cite this version:

Mélanie Blazère, Fabrice Gamboa, Jean-Michel Loubes. PLS: a new statistical insight through the prism of orthogonal polynomials. 2014. <hal-00995049>

HAL Id: hal-00995049

<https://hal.archives-ouvertes.fr/hal-00995049>

Submitted on 22 May 2014

HAL is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.

Partial Least Square

A new statistical insight through the prism of orthogonal polynomials

Mélanie Blazère, Jean-Michel Loubes and Fabrice Gamboa

Abstract

Partial Least Square (PLS) is a dimension reduction method used to remove multicollinearities in a regression model. However contrary to Principal Components Analysis (PCA) the PLS components are also chosen to be optimal for predicting the response Y . In this paper we provide a new and explicit formula for the residuals. We show that the residuals are completely determined by the spectrum of the design matrix and by the noise on the observations. Because few are known on the behaviour of the PLS components we also investigate their statistical properties in a regression context. New results on regression and prediction error for PLS are stated under the assumption of a low variance of the noise.

Keywords

Partial Least Square, multivariate regression, multicollinearity, dimension reduction, constrained least square, orthogonal polynomials, prediction error.

1 Introduction

Partial Least Square (PLS), introduced in 1985 by [Wold \(1985\)](#), is nowadays a widely used dimension reduction technique in multivariate analysis especially when we have to handle high dimensional or highly correlated data in a regression context. Originally designed to remove the problem of multicollinearity in the set of explanatory variables, PLS acts as a dimension reduction method by creating a new subset of variables which are also optimal for predicting the output variable. During the last decades this method has been developed and studied for a large part by [Helland \(1988\)](#) and by [Frank and Friedman \(1993\)](#). Partial Least Square was originally developed for chemometrics applications (see for example [Wold et al. \(2001\)](#) and [Frank and Friedman \(1993\)](#)) but gained attention in biosciences in particular in the analysis of high dimensional genomic data. We refer for instance to [Boulesteix and Strimmer \(2007\)](#) or to [Lê Cao et al. \(2008\)](#) for various applications in this field.

If the PLS method proved helpful in a large variety of situations, this iterative procedure is complex and little is known about its theoretical properties but PLS has been well investigated by practical experiments. To name just a few, [Naes and Martens \(1985\)](#) discussed theoretical and computational considerations of PLS and PCR (Principal Component Regression) on simulated and real data. [Frank and Friedman \(1993\)](#) provided a heuristic comparison of the performances of OLS, PCR, Ridge regression and PLS in different situations. [Garthwaite \(1994\)](#) compared PLS with four other methods (ordinary least

squares, forward variable selection, principal components regression, and a Stein shrinkage method) through simulations. Only very recently, some theoretical insights have been given by [Delaigle and Hall \(2012\)](#) for functional data.

In this work, we provide a new direction to analyze some statistical aspects of the PLS method. For this, we will draw connections between PLS, Krylov subspaces and the regularization of inverse problems (see [Engl et al. \(1996\)](#)).

The paper falls into the following sections. In [Section 2](#) we present the framework within which we study PLS and we briefly recall what is the PLS method. We also highlight the connection between PLS and Krylov subspaces. Because the directions of the new subspace onto which we project the observations depend on the response variable it is quite difficult to study the statistical properties of PLS using just the algorithmic construction of the new subspace. In this paper we adopt the point of view of PLS viewed as a constrained least square problem and use its connection with inverse problem with a statistical point of view. In [Section 3](#) we highlight the connection between PLS and the minimization of the least squares over polynomial subspaces. Then in [Section 4](#) we provide a new formulation of the residuals for each direction defined by the eigenvectors of the covariance matrix. The interest of such a formulation rests on the fact that it provides an explicit expression of the residuals in terms of both the noise on the observations and on the eigenelements of the covariance matrix. This formulation will enable a further study of the PLS method performance in a regression framework. In [Section 5](#) we study PLS in the context of a high-dimensional multiple regression model. We first define the model under study. Then we detail our main results for noisy sample and new statistical aspects of PLS. We provide bounds for the empirical risk and for the mean square error of prediction under the assumption of a low variance of the noise. Asymptotic properties of the prediction error are also discussed. We also highlight the limitations of this method according to the features of the data.

2 Presentation of the framework

2.1 Notation and remarks

We first introduce some of the notation we use in this paper. By $\langle x, y \rangle$ we denote the inner product between the vectors $x, y \in \mathbb{R}^n$. The transpose of a matrix A is denoted by A^T and it depends on the underlying inner product, i.e. $\langle Ax, y \rangle = \langle x, A^T y \rangle$. The induced vector norm is $\|x\| = \sqrt{\langle x, x \rangle}$. In most cases we work with the Euclidean inner product i.e. $\langle x, y \rangle = x^T y$ and the induced norm is the ℓ_2 -norm. For any positive definite matrix M , the M -inner product is defined as $\langle x, y \rangle_M = x^T M y$ and the operator norm is given by $\|M\| = \max_{\|x\|=1} \|Mx\|$. We simply denote by I the identity matrix with the corresponding dimension. For every $k \in \mathbb{N}$ we denote by \mathcal{P}_k the set of the polynomials of degree less than k and by $\mathcal{P}_{k,1}$ the set of the polynomial of degree less than k whose constant term equals 1.

The figures which appear in this paper are the result of simulations which have been performed with R using the package `pls.genomics` developed by Boulesteix and al. The function `pls.regression` has been used to fit the model.

2.2 The regression model

We consider the following regression model

$$Y = X\beta^* + \varepsilon \tag{1}$$

where

- $Y = (Y_1, \dots, Y_n)^T \in \mathbb{R}^n$ is the vector of the observed outcome, also called the response.
- $X = (X_{ij})_{1 \leq i \leq n, 1 \leq j \leq p} \in \mathbb{M}_{n \times p}$ is the design matrix which is considered as fixed and contains the predictors.
- $\beta^* = (\beta_1^*, \dots, \beta_p^*)^T \in \mathbb{R}^p$ is the unknown parameter vector and represents the variables of interest.
- $\varepsilon = (\varepsilon_1, \dots, \varepsilon_n)^T \in \mathbb{R}^n$ captures the noise.

In other word we are concerned with finding a good approximation $\hat{\beta}$ of the solution β^* of the above linear problem where only noisy observations are available. For the moment we only assume that the real variables $\varepsilon_1, \dots, \varepsilon_n$ are unobservable i.i.d random variables. We allow p to be much larger than n i.e $p \gg n$. We denote by r the rank of $X^T X$. Of course $r \leq \min(n, p)$. The aim is to estimate the unknown parameter β^* from the observations of the pairs $(Y_i, X_i)_{1 \leq i \leq n}$. The usual ordinary least squares (OLS) estimates β^* by $\hat{\beta}$ where

$$\hat{\beta} \in \underset{\beta \in \mathbb{R}^p}{\operatorname{argmin}} \|Y - X\beta\|^2.$$

However we know that in the case of highly correlated explanatory variables and/or when the explanatory variables outnumber the observations i.e $p \gg n$ the regression model is ill conditioned and the OLS estimator behaves badly. The estimated parameter can be very unstable and far from the target leading to unaccurate predictions. To remove the problem of multicollinearities in regression model a solution consists of creating latent variables using Principal Component Analysis (PCA). However, the new variables are chosen to explain X but they may not explain Y well. [Jolliffe \(1982\)](#) provided several real-life examples where the principal components corresponding to small eigenvalues have high correlation with Y . To avoid this problem a possible solution is to use Partial Least Square which has been heavily promoted as an alternative to OLS in the literature.

2.3 The PLS method

In this subsection we briefly recall the method. The PLS method, introduced by [Wold \(1985\)](#), emerged in order to remove the problem of multicollinearity in a regression model (when the number of covariates is large or when there are dependancies between them). In fact PLS is a statistical method whose challenge is to find principal components that explain X as well as possible and are also good predictors for Y .

The PLS method at step K (where $K \leq r$) consists in finding $(w_k)_{1 \leq k \leq K}$ and $(t_k)_{1 \leq k \leq K}$ which maximize $[\operatorname{Cov}(Y, Xw_k)]^2$ under the constraint

1. $\|w_k\|^2 = w_k^T w_k = 1$
2. $t_k = Xw_k$ is orthogonal to t_1, \dots, t_{k-1} i.e $w_k^T X^T X w_l = 0$ for $l = 1, \dots, k-1$.

Therefore the PLS method is a procedure which iteratively constructs a subspace of dimension K (spanned by $(w_k)_{1 \leq k \leq K}$) in such a way that the new latent variables $(t_k)_{1 \leq k \leq K}$ (which are the projections of the original ones) maximize both the correlation with the response and the variance of the explanatory variables. The original algorithms were developed by [Wold et al. \(1983\)](#) and a decade later by [Martens and Naes \(1992\)](#).

Once the latent variables $(t_k)_{1 \leq k \leq K}$ are built, one can compute the linear regression of Y on t_1, \dots, t_K to estimate β^* . We can notice that this method is of particular interest because it can analyze data with strongly correlated, noisy and numerous covariates. Furthermore dimension reduction and regression are performed simultaneously. We refer to [Helland \(1990, 1988, 2001\)](#) for the study of the main properties of PLS and to [Krämer \(2007\)](#) for a complete overview on the recent advances on PLS. Proposition 2.1 below recalls what appears to us as one of the main result on PLS because it is the starting point of our work. This proposition shows that the PLS estimator at step K is defined as the argument which minimizes the least square over some particular subspace of dimension K .

Proposition 2.1. [Helland \(1990\)](#)

$$\hat{\beta}_K^{PLS} = \underset{\beta \in \mathcal{K}^K(X^T X, X^T Y)}{\operatorname{argmin}} \|Y - X\beta\|^2 \quad (2)$$

where $\mathcal{K}^k(X^T X, X^T Y) = \{X^T Y, (X^T X)X^T Y, \dots, (X^T X)^{k-1}X^T Y\}$.

The space $\mathcal{K}^k(X^T X, X^T Y)$ spanned by $X^T Y, (X^T X)X^T Y, \dots, (X^T X)^{k-1}X^T Y$ (and denoted by \mathcal{K}^k when there is no possible confusion) is called the k^{th} Krylov subspace. We refer to [Saad \(1992\)](#) for a further study of these spaces. We can notice that, as for PCR, PLS is a constrained least square estimator where the constraints are not on the norm of the parameter (as for Ridge regression or for the Lasso) but are linear constraints which ensure that the estimated parameter belongs to the Krylov subspace associated to $X^T X$ and to $X^T Y$. However we have to be careful that contrary to PCR the PLS linear constraints are random.

Using this connection with Krylov subspaces [Phatak and de Hoog \(2002\)](#) showed that the PLS iterates are the same as the ones of the Conjugate Gradient(CG). Thus PLS can also be viewed as CG applied in the statistical framework of linear regression models. Phatak and de Hoog also used the connection between CG, Lanczos method and PLS to give simpler proofs of two known results. The first one is the shrinkage properties of PLS ($\|\hat{\beta}_k^{PLS}\| \leq \|\hat{\beta}_{k+1}^{PLS}\|$) proved by [De Jong \(1995\)](#) and the second is the fact that PLS fits better than PCR ($\|\hat{\beta}_k^{PLS}\| \leq \|\hat{\beta}_{OLS}\|$) proved by [Goutis \(1996\)](#).

3 A close connection between PLS and orthogonal polynomials

In this section we show that the PLS solution can be written as the polynomial solution of a minimization problem. Then we prove that the sequence of the residuals in each eigenvectors direction can be expressed as a sequence of orthogonal polynomials with respect to a discrete measure. This measure depends on the eigenvalues of the design matrix and on the projection of the response onto the associated eigenvectors.

3.1 A useful tool to analyze the properties of PLS

Consider the singular value decomposition of X given by

$$X = UDV^T$$

where

- $U^T U = I$ and u_1, \dots, u_p are the columns of U and form an orthonormal basis of \mathbb{R}^n .

- $V^T V = I$ and v_1, \dots, v_p are the columns of V and form an orthonormal basis of \mathbb{R}^n .
- $D \in \mathbb{M}_{n,p}$ is a matrix which contains $(\sqrt{\lambda_1}, \dots, \sqrt{\lambda_n})$ on the diagonal and zero anywhere else.
- We assume that $\lambda_1 \geq \lambda_2 \geq \dots \geq \lambda_n > 0 = \lambda_{n+1} = \dots = \lambda_p$. In other words we assume that $X^T X$ is of full rank i.e of rank n .

Of course we have $X^T u_i = \sqrt{\lambda_i} v_i$, $XX^T u_i = \lambda_i u_i$ for all $i = 1, \dots, n$ and $Xv_i = \sqrt{\lambda_i} u_i$, $X^T X v_i = \lambda_i v_i$ for all $i = 1, \dots, p$. We define $\tilde{\varepsilon}_i := \varepsilon^T u_i$, $i = 1, \dots, n$ and $\tilde{\beta}_i^* := \beta^{*T} v_i$, $i = 1, \dots, p$ the projections of ε and β^* respectively onto the right and left eigenvectors of X .

We assume that $k \leq \mu + 1$ where μ is the grade of $X^T Y$ with respect to $X^T X$ i.e the degree of the nonzero monic polynomial P of lowest degree such that $P(X^T X)X^T Y = 0$. The Krylov subspace \mathcal{K}^k is of dimension k if and only if $\mu \geq k - 1$ in such a way that in this case we have $\dim(\mathcal{K}^k) = k$.

We can notice that the maximal dimension of the Krylov subspace sequence is also linked to the number of non zero eigenvalues λ_i for which $u_i^T Y \neq 0$ (see [Helland \(1990\)](#)). These particular eigenvalues are called the relevant eigenvalues. If the number of relevant eigenvalues is n then the maximal dimension of the Krylov subspaces sequence is also n and for all $k \leq n$ the dimension of \mathcal{K}^k is exactly k . In particular if $X^T Y = \sum_{i=1}^k \sqrt{\lambda_i} (u_i^T Y) v_i$ then the PLS iterations will terminate in at most k iterations.

3.2 Link with the regularization of inverse problems methods: a minimization problem over polynomials

We recall that \mathcal{P}_k is the set of the polynomials of degree less than k and $\mathcal{P}_{k,1}$ the set of the polynomial of degree less than k whose constant term equals one. By combining formula (2) which expresses the PLS estimator as a constrained least square over Krylov subspace with the definition of \mathcal{K}^k it is easy to show that the PLS estimator can also be expressed as the solution of a minimization problem over polynomials.

Proposition 3.1. *For $k \leq n$ we have*

$$\hat{\beta}_k = \hat{P}_k(X^T X)X^T Y \quad (3)$$

where \hat{P}_k is a polynomial of degree less than $k - 1$ which satisfies

$$\|Y - X\hat{P}_k(X^T X)X^T Y\|^2 = \underset{P \in \mathcal{P}_{k-1}}{\operatorname{argmin}} \|Y - XP(X^T X)X^T Y\|^2$$

and

$$\|Y - X\hat{\beta}_k\|^2 = \|\hat{Q}_k(XX^T)Y\|^2 = \underset{Q \in \mathcal{P}_{k,1}}{\min} \|Q(XX^T)Y\|^2 \quad (4)$$

where $\hat{Q}_k(t) = 1 - t\hat{P}_k(t)$ is a polynomial of degree less than k and of constant term equals to one.

Notice that for all $k \leq n$ we have $X\hat{\beta}_k = \hat{\Pi}_k Y$ where $\hat{\Pi}_k$ is the orthogonal projector onto the random space $\mathcal{K}^k(XX^T, XX^T Y)$ of dimension k . In particular for $k = n$, we have $X\hat{\beta}_n = Y$ and $\|Y - X\hat{\beta}_n\|^2 = 0$ because $\mathcal{K}^n(XX^T, XX^T Y)$ is of dimension n . In the following we will omit this trivial case.

Proposition 3.1 shows that the PLS method is another regularization method for ill-posed inverse problem (see [Engl et al. \(1996\)](#)). In fact when the explanatory variables

are highly correlated or when they outnumber the number of observations the regression model we consider is ill-posed. The idea behind PLS is to approximate the ill-posed problem by a family of nearby well-posed problem by seeking for regularization operator \mathcal{R}_α such that $\mathcal{R}_\alpha(X^T X)X^T Y$ close to β^* where \mathcal{R}_α is under a polynomial form. In fact the polynomials \hat{P}_k play the role of \mathcal{R}_α and $\alpha = k$ is the regularization parameter. We also refer to [Blanchard and Mathé \(2012\)](#) to take a more in depth look on statistical inverse problems and Conjugate Gradient because PLS is closely related to Conjugate Gradient with a statistical point of view.

The idea of considering the Krylov subspace and thus polynomial approximation is at the heart of the issue for PLS. We present below a result which gives a good reason to search for polynomial approximations. Indeed the theorem of Cayley-Hamilton tells us that we can represent the inverse of a nonsingular matrix A in terms of the powers of A . It is no longer the case for a singular matrix because the inverse does not exist. But the idea behind PLS remains quite the same for non singular matrix. It consists of using Krylov subspaces to approximate the pseudo inverse as a polynomial in A . In fact according to (3) the PLS estimator $\hat{\beta}_k$ is of the form $\hat{P}_k(X^T X)X^T Y$ where \hat{P}_k is a polynomial of degree less than $k - 1$ and thus consists in a kind of regularization of the inverse of $X^T X$. Notice that since $\dim(\mathcal{K}^k) = k$ the polynomial \hat{P}_k is in fact of degree exactly $k - 1$. If $X^T X$ is invertible then the PLS method generates a sequence of polynomial approximation of the inverse of $X^T X$ and when $k = n$ we recover the inverse of $X^T X$ exactly.

So PLS is also equivalent to finding an optimal polynomial \hat{Q}_k of degree k with $\hat{Q}_k(0) = 1$ minimizing $\|Q(XX^T)Y\|^2$. Notice that if there exists a polynomial Q of degree k with $Q(0) = 1$ small on the spectrum of XX^T then $\|Y - X\hat{\beta}_k\|^2$ will be small too. In particular if the eigenvalues are clustered into k groups (i.e can be divided into k groups whose diameter are very small) then $\|Y - X\hat{\beta}_k\|^2$ has a good chance to be small as well. The polynomial \hat{Q}_k quantifies the quality of the approximation of the response Y at the k^{th} step. We call these polynomials the residuals.

3.3 Link with orthogonal polynomials

In this subsection we first prove that the sequence of polynomials $(\hat{Q}_k)_{1 \leq k \leq n}$ defined in Proposition 3.1 is orthogonal with respect to a discrete measure denoted by $\hat{\mu}$.

Proposition 3.2. *$\hat{Q}_1, \hat{Q}_2, \dots, \hat{Q}_{n-1}$ is a sequence of orthonormal polynomials with respect to the measure*

$$d\hat{\mu}(\lambda) = \sum_{j=1}^n \lambda_j (u_j^T Y)^2 \delta_{\lambda_j}.$$

The support of the measure $\hat{\mu}$ consists of the $(\lambda_i)_{1 \leq i \leq n}$ and the weights depend on $(\lambda_i)_{1 \leq i \leq n}$ and $(u_j^T Y)_{1 \leq j \leq n}$. These last quantities capture both the variation in X and the correlation between X and Y .

4 A new expression for the residuals in the eigenvectors direction

4.1 Main Result

If the PLS properties are not completely understood it is partly because the solution is a non linear function of the data Y . PLS is an iterative method and therefore if we

perturb Y the perturbation propagates through the sequence of Krylov subspaces in a non linear way which makes difficult the explicit study of the PLS estimator. In this section we provide a new explicit and exact formulation of the residuals which clearly shows how the disturbance on the observations impacts on the residuals.

In this section a new and exact expression for $\hat{Q}_k(\lambda_i)$ is proposed for all $k = 1, \dots, n-1$ and all $i = 1, \dots, n$.

Theorem 4.1. *Let $k \leq n$ and*

$$I_k^+ = \{n \geq j_1 > \dots > j_k \geq 1\}.$$

We have

$$\hat{Q}_k(\lambda_i) = \sum_{(j_1, \dots, j_k) \in I_k^+} \left[\frac{\hat{p}_{j_1}^2 \dots \hat{p}_{j_k}^2 \lambda_{j_1}^2 \dots \lambda_{j_k}^2 V(\lambda_{j_1}, \dots, \lambda_{j_k})^2}{\sum_{(j_1, \dots, j_k) \in I_k^+} \hat{p}_{j_1}^2 \dots \hat{p}_{j_k}^2 \lambda_{j_1}^2 \dots \lambda_{j_k}^2 V(\lambda_{j_1}, \dots, \lambda_{j_k})^2} \right] \prod_{l=1}^k \left(1 - \frac{\lambda_i}{\lambda_{j_l}}\right). \quad (5)$$

where $\hat{p}_i := p_i + \tilde{\varepsilon}_i$ with $p_i := (X\beta^*)^T u_i = \sqrt{\lambda_i} \tilde{\beta}_i^*$ and $\tilde{\varepsilon}_i := \varepsilon^T u_i$.

For all $k < n$ we recover that \hat{Q}_k is a polynomial of degree k and $\hat{Q}_k(0) = 1$. The expression of $\hat{Q}_k(\lambda_i)$ given in Proposition 7.2 depends explicitly on the observations noise and on the eigenelements of X contrary to the expression provided in the paper of [Lingjaerde and Christophersen \(2000\)](#). Formula (5) is also valid for $k = n$ but in this case we recover that $\hat{Q}_n(\lambda_i) = 0$ for all $i = 1, \dots, n$.

Now assume that there are only k distinct eigenvalues among the n ones and denote by $\tilde{\lambda}_1, \dots, \tilde{\lambda}_k$ the different representatives. Then for all $i = 1, \dots, n$ formula (5) implies

$$\hat{Q}_k(\lambda_i) = \prod_{j=1}^k \left(1 - \frac{\lambda_i}{\tilde{\lambda}_j}\right) = 0.$$

Thus the residuals along each eigenvectors equal zero at step k if there are less than k different non zero eigenvalues (notice that this is of course the case when $k = n$). Furthermore if we assume that there exists only k eigenvectors denoted by $u_{j_1}^-, \dots, u_{j_k}^-$ such that $\hat{p}_{j_1}^- \neq 0, \dots, \hat{p}_{j_k}^- \neq 0$ then formula (5) becomes

$$\hat{Q}_k(\lambda_i) = \prod_{j=1}^k \left(1 - \frac{\lambda_i}{\lambda_{j_l}^-}\right).$$

Therefore for all $\lambda \in \{\lambda_{j_1}^-, \dots, \lambda_{j_k}^-\}$, $\hat{Q}_k(\lambda) = 0$ and thus we find $\|Y - X\hat{\beta}_k\|^2 = \|\hat{Q}_k(XX^T)Y\|^2 = \sum_{i=1}^n \hat{Q}_k(\lambda_i)^2 (u_i^T Y)^2 = 0$.

For all $((j_1, \dots, j_k)) \in I_k^+$, let

$$\hat{w}_{j_1, \dots, j_k} := \frac{\hat{p}_{j_1}^2 \dots \hat{p}_{j_k}^2 \lambda_{j_1}^2 \dots \lambda_{j_k}^2 V(\lambda_{j_1}, \dots, \lambda_{j_k})^2}{\sum_{(j_1, \dots, j_k) \in I_k^+} \hat{p}_{j_1}^2 \dots \hat{p}_{j_k}^2 \lambda_{j_1}^2 \dots \lambda_{j_k}^2 V(\lambda_{j_1}, \dots, \lambda_{j_k})^2}$$

and

$$g_{j_1, \dots, j_k}(x) = \prod_{l=1}^k \left(1 - \frac{x}{\lambda_{j_l}}\right).$$

Notice that this last function is again a polynomial in x of degree k whose constant term is equal to one and is zero at $\lambda_{j_1}, \dots, \lambda_{j_k}$ which are elements in the spectrum of XX^T . We have

$$\hat{Q}_k(\lambda_i) = \sum_{(j_1, \dots, j_k) \in I_k^+} \hat{w}_{(j_1, \dots, j_k)} g_{j_1, \dots, j_k}(\lambda_i).$$

Besides, for all $(j_1, \dots, j_k) \in I_k^+$, $0 < \hat{w}_{(j_1, \dots, j_k)} \leq 1$ and $\sum_{(j_1, \dots, j_k) \in I_k^+} \hat{w}_{(j_1, \dots, j_k)} = 1$. Thus the weights $(\hat{w}_{(j_1, \dots, j_k)})_{I_k^+}$ are probabilities. Therefore $\hat{Q}_k(\lambda_i)$ is the sum over all elements in I_k^+ of $g_{j_1, \dots, j_k}(\lambda_i)$ weighted by the probabilities $\hat{w}_{(j_1, \dots, j_k)}$. It is a kind of barycenter of all the polynomials in $\mathcal{P}_{k,1}$ whose roots are subsets of $\{\lambda_1, \dots, \lambda_n\}$. The weights are not easy to interpret but they are even greater when the magnitude and the distance between the involved eigenvalues are large and the contribution of the response along the associated eigenvectors is important. From formula (5) we can state in a very large way that

$$|\hat{Q}_k(\lambda_i)| \leq \max_{I_k^+} \left(\prod_{l=1}^k \left| 1 - \frac{\lambda_i}{\lambda_{j_l}} \right| \right).$$

In particular if $\lambda_1(1 - \varepsilon) \leq \lambda_i \leq \lambda_n(1 + \varepsilon)$ then

$$|\hat{Q}_k(\lambda_i)| \leq \varepsilon^k.$$

Here is an example of the residuals path with respect to the eigenvectors directions for 100 nonzero eigenvalues which are distributed around 10 different values. We also represent the residuals only for the extremal eigenvalues to better see the difference of behaviour.

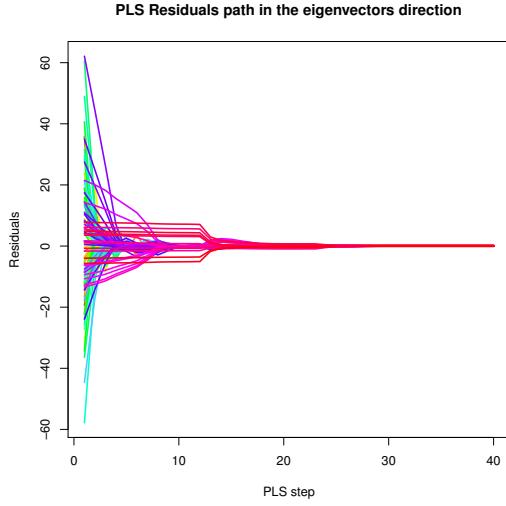


Figure 1

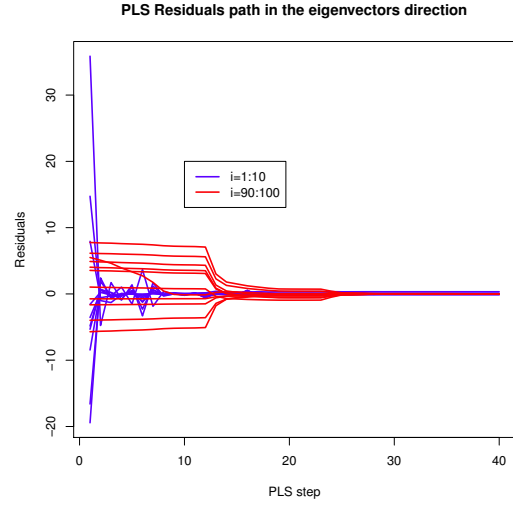


Figure 2

Proposition 4.2. *Let n and k fixed and $i \in \llbracket 1, n \rrbracket$. If $\lambda_j = \lambda_i + \delta$ then*

$$|\hat{Q}_k(\lambda_i) - \hat{Q}_k(\lambda_j)| \leq \delta \max_{I_k^+} \left[\sum_{l=1}^k \frac{1}{\lambda_{j_l}} \prod_{m \neq l} \left(1 - \frac{\lambda_i}{\lambda_{j_m}} \right) \right] + O(\delta^2).$$

Proof. Let n and k fixed and $i \in \llbracket 1, n \rrbracket$. Assume that $\lambda_j = \lambda_i + \delta$. We have

$$\hat{Q}_k(\lambda_j) = \sum_{(j_1, \dots, j_k) \in I_k^+} \left[\hat{w}_{(j_1, \dots, j_k)} \prod_{l=1}^k \left(1 - \frac{\lambda_j}{\lambda_{j_l}} \right) \right] = \sum_{(j_1, \dots, j_k) \in I_k^+} \left[\hat{w}_{(j_1, \dots, j_k)} \prod_{l=1}^k \left(1 - \frac{\lambda_i + \delta}{\lambda_{j_l}} \right) \right].$$

By expanding $\prod_{l=1}^k \left(1 - \frac{\lambda_i + \delta}{\lambda_{j_l}} \right)$ we get

$$\hat{Q}_k(\lambda_i) = \hat{Q}_k(\lambda_j) - \delta \sum_{(j_1, \dots, j_k) \in I_k^+} \hat{w}_{(j_1, \dots, j_k)} \left[\sum_{l=1}^k \frac{1}{\lambda_{j_k}} \prod_{m \neq l} \left(1 - \frac{\lambda_i}{\lambda_{j_m}} \right) \right] + O(\delta^2).$$

Then using the fact that $\sum_{(j_1, \dots, j_k) \in I_k^+} \hat{w}_{(j_1, \dots, j_k)} = 1$ we deduce Proposition 4.2. \square

Thus for nearby eigenvalues the filter factors are almost the same. Therefore if $\hat{Q}_k(\lambda_i)$ is small then $\hat{Q}_k(\lambda_j)$ will be small too if λ_j is closed enough to λ_i . In particular if the eigenvalues are clustered into k groups and if the residuals associated to the center of the clusters are close to zero then all the residuals will be closed to zero too.

The expression of the residuals provided by Theorem 4.1 will be very useful and central elsewhere in this paper to further explore the PLS method and prove new statistical results.

4.2 Filter factors and shrinkage properties

In this subsection we show that we recover some of the results first proved by [Butler and Denham \(2000\)](#) and [Lingjaerde and Christophersen \(2000\)](#) on the shrinkage properties of the PLS estimator and more particularly on its expansion or contraction in the eigenvectors directions using the expression of the residuals provided by Theorem 4.1. We have (see [Lingjaerde and Christophersen \(2000\)](#))

$$\hat{\beta}_k = \sum_{i=1}^n f_i^k \frac{\hat{p}_i}{\sqrt{\lambda_i}} v_i$$

where the elements $f_i^k = 1 - \hat{Q}_k(\lambda_i)$ are called the filter factors. In their study Lingjaerde and Christophersen use the following implicit expression of \hat{Q}_k

$$\hat{Q}_k(t) = \frac{(\theta_1^{(k)} - t) \dots (\theta_k^{(k)} - t)}{\theta_1^{(k)} \dots \theta_k^{(k)}}$$

where $(\theta_i^{(k)})_{1 \leq i \leq n}$ are the eigenvalues of $W_k(W_k^T \Sigma W_k)W_k^T$ (also called the Ritz eigenvalues) to study the shrinkage properties of PLS. [Lingjaerde and Christophersen \(2000\)](#) showed that all the PLS shrinkage factors are not in $[0, 1]$ and can be larger than one. They even proved more precisely that the filter factors oscillate between below and above one (depending on the parity of the index of the factors). We recover these results from our expression of the residuals provided in Theorem 4.1

$$\hat{Q}_k(\lambda_i) = \sum_{(j_1, \dots, j_k) \in I_k^+} \left[\hat{w}_{(j_1, \dots, j_k)} \prod_{l=1}^k \left(1 - \frac{\lambda_i}{\lambda_{j_l}} \right) \right],$$

where we recall that $\hat{w}_{(j_1, \dots, j_k)} := \frac{\hat{p}_{j_1}^2 \dots \hat{p}_{j_k}^2 \lambda_{j_1}^2 \dots \lambda_{j_k}^2 V(\lambda_{j_1}, \dots, \lambda_{j_k})^2}{\sum_{(j_1, \dots, j_k) \in I_k^+} \hat{p}_{j_1}^2 \dots \hat{p}_{j_k}^2 \lambda_{j_1}^2 \dots \lambda_{j_k}^2 V(\lambda_{j_1}, \dots, \lambda_{j_k})^2}$. From this formula we deduce

$$f_i^k = \sum_{(j_1, \dots, j_k) \in I_k^+} \hat{w}_{(j_1, \dots, j_k)} \left[1 - \prod_{l=1}^k \left(1 - \frac{\lambda_i}{\lambda_{j_l}} \right) \right].$$

Notice that the filter factors are completely and explicitly determined by the spectrum and the eigenvectors of $X^T X$.

If $k < n$ and $i = n$ then $0 < \prod_{l=1}^k (1 - \frac{\lambda_n}{\lambda_{j_l}}) < 1$ and from

$$\sum_{(j_1, \dots, j_k) \in I_k^+} \hat{w}_{(j_1, \dots, j_k)} = 1$$

we conclude that $0 < f_n^k < 1$.

If $k < n$ and $i = 1$ then

$$\begin{cases} \prod_{l=1}^k (1 - \frac{\lambda_1}{\lambda_{j_l}}) < 0 & \text{if } k \text{ is odd} \\ \prod_{l=1}^k (1 - \frac{\lambda_1}{\lambda_{j_l}}) > 0 & \text{if } k \text{ is even} \end{cases} \quad (6)$$

and

$$\begin{cases} f_1^k > 1 & \text{if } k \text{ is odd} \\ f_1^k < 1 & \text{if } k \text{ is even.} \end{cases} \quad (7)$$

For the other filter factors we can have $f_i^k \leq 1$ or $f_i^k \geq 1$ (depending on the distribution of the spectrum) contrary to the PCR or Ridge filter factors which always lies in $[0, 1]$. Therefore PLS shrinks in some directions and expands in others. However the PLS estimator is considered as a shrinkage estimator because $\|\hat{\beta}_k^{PLS}\| \leq \|\hat{\beta}_{OLS}\|$ (see [Goutis \(1996\)](#)).

We also recover Theorem 7 of [Lingjaerde and Christoffersen \(2000\)](#). Indeed if we have $\lambda_i < \lambda_n(1 + \sqrt{\epsilon})$ then a straightforward calculation using formula (5) leads to $f_i^k < 1 + \epsilon$.

5 Bounds for the empirical risk and prediction error

In this section we further explore the statistical properties of PLS. For this, we investigate the accuracy of PLS through the study of the empirical risk and the least square error of prediction which are two criteria commonly used for assessing the quality of an estimator.

From now, on we assume that the $(\varepsilon_i)_{1 \leq i \leq n}$ are i.i.d centered random variables with common variance σ^2 and for simplicity we also assume that the observations on the X variables are centered and normalized, that is $\frac{1}{n} \sum_{i=1}^n X_{ij} = 0$ and $\frac{1}{n} \sum_{i=1}^n X_{ij}^2 = 1$.

5.1 Empirical risk

The following proposition provides an upper bound for the MSE (mean square error). The MSE quantifies the fit of the model to the data set used.

Proposition 5.1. *We have for $k < n$*

$$\mathbb{E} \left[\frac{1}{n} \|Y - X\hat{\beta}_k\|^2 \right] \leq \left[\frac{1}{n} \left(\frac{\sqrt{C(X^T X)} - 1}{\sqrt{C(X^T X)} + 1} \right)^{2k} \|X\beta^*\|^2 \right] \left[1 + \frac{n\sigma^2}{\|X\beta^*\|^2} \right] \quad (8)$$

where $C(X^T X) = \frac{\lambda_1}{\lambda_n}$ is the ratio of the two extreme non zero eigenvalues of $X^T X$.

Obviously, if $k = n$,

$$\mathbb{E} \left[\frac{1}{n} \|Y - X\hat{\beta}_n\|^2 \right] = 0.$$

The first factor in equation (8) represents the error due only to the regularization if no noise (projection of $X\beta^*$ onto the Krylov subspace $\mathcal{K}^k(X^T X, X^T X\beta^*)$ and in the second factor $\frac{n\sigma^2}{\|X\beta^*\|^2}$ represents the inverse of the signal to noise ratio. We can notice that the upper bound relies on $\|X\beta^*\|^2 = \sum_{i=1}^n \lambda_i \tilde{\beta}_i^2$. This term links the regularity of β^* with the decay of the eigenvalues of $X^T X$. It thus can be seen as a Source Condition, see for instance in Engl et al. (1996). We can state a result similar to the one of Proposition 5.1 replacing $\|\cdot\|_2$ by $\|\cdot\|_p$, $p \in \mathbb{N}^*$.

We can notice that the convergence of the empirical risk is associated with upper bounds derived using scaled and shifted Chebyshev polynomials. In fact the key of the proof of Proposition 5.1 is essentially based on the following proposition

Proposition 5.2. *Saad (1992)*

Let $[\alpha, \beta]$ be a non empty interval in \mathbb{R} and let γ be any scalar such with $\gamma \notin]\alpha, \beta[$. We define $\mathcal{E}_k := \{P \text{ polynomial of degree } k \text{ with } P(\gamma) = 1\}$.

Then the minimum $\min_{P \in \mathcal{E}_k} \max_{t \in [\alpha, \beta]} |P(t)|$ is reached by the polynomial

$$\hat{C}_k(t) = \frac{C_k(1 + 2\frac{t-\beta}{\beta-\alpha})}{C_k(1 + 2\frac{\gamma-\beta}{\beta-\alpha})},$$

where C_k is the k^{th} Chebychev polynomial i.e., for $x \in [-1, 1]$,

$$C_k(x) = \frac{1}{2} \left((x - \sqrt{x^2 - 1})^k + (x + \sqrt{x^2 - 1})^k \right).$$

The maximum of C_k for $x \in [-1, 1]$ is 1 and

$$\min_{P \in \mathcal{E}_k} \max_{t \in [\alpha, \beta]} |P(t)| = \frac{1}{|C_k(1 + 2\frac{\gamma-\beta}{\beta-\alpha})|} = \frac{1}{|C_k(2\frac{\gamma-\mu}{\beta-\alpha})|}$$

with $\mu = \frac{\alpha+\beta}{2}$.

Notice that the convergence rate of the empirical risk is exponential in k with respect to the ratio between the maximum and the minimum of the non zero eigenvalues. In fact $\left| \frac{\sqrt{C(X^T X)-1}}{\sqrt{C(X^T X)+1}} \right| < 1$ and is equal to zero if and only if all the eigenvalues are the same. Therefore the closer to one is the condition number $C(X^T X)$ the faster is the decrease of the empirical risk with respect to k , in an exponential way while it turns polynomial for Ridge Regression for instance.

Figure 3 represents the empirical risk for different values of the level of noise.

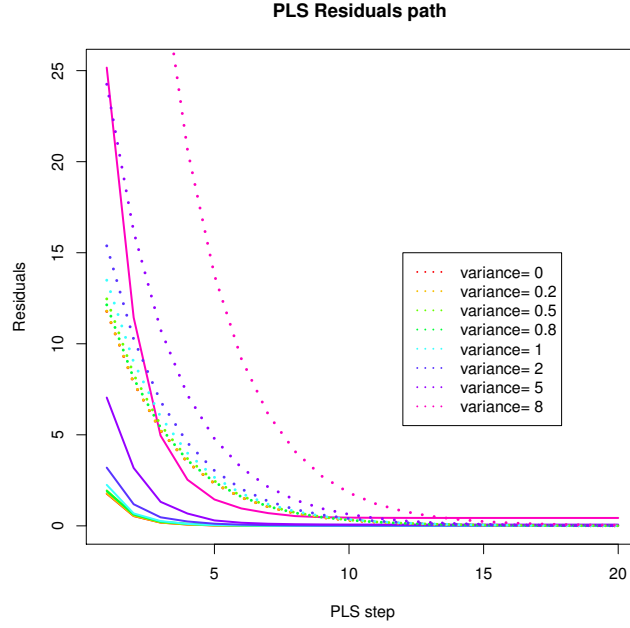


Figure 3

A way to improve the method could be the use of preconditioners i.e. of a matrix M used to convert the problem $X^T Y = X^T X \beta^* + X^T \varepsilon$ into another equivalent problem i.e. into $M^{-1} X^T X \beta = M^{-1} X^T (Y - \varepsilon)$ in such a way that it increases the rate of convergence.

As noticed below if there are only k distinct eigenvalues or if the contribution of Y is only non zero along k eigenvectors then we also have $\mathbb{E} \left[\frac{1}{n} \|Y - X \hat{\beta}_k\|^2 \right] = 0$. This is a straightforward consequence of equation (4) in Proposition 3.1, taking $Q(x) = \prod_{i=1}^k \left(1 - \frac{x}{\lambda_i} \right)$ where $(\bar{\lambda}_i)_{1 \leq i \leq k}$ are the representatives of respectively the different non zero eigenvalues and the eigenvalues associated to a non zero contribution of the response to the associated eigenvectors. In the same way the empirical risk will be very small at step k if the eigenvalues are clustered into k groups.

To illustrate this particular behaviour of the PLS estimator we have performed some simulations. The data sets are simulated according to model (1) with $n = p = 100$. The best latent components are chosen using the function `pls.regression.cv`. For the first simulation below we consider that the eigenvalues are partitionned into 2 clusters and for the second one that the eigenvalues are partitionned into 10 clusters.

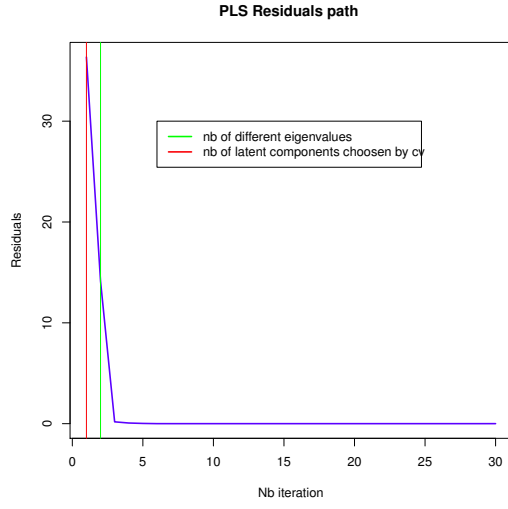


Figure 4

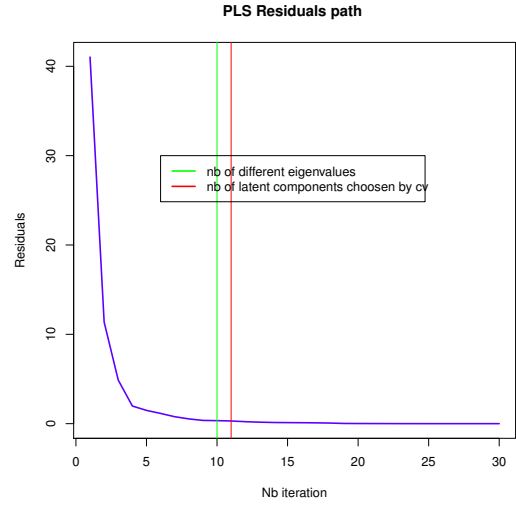


Figure 5

For the next three models we consider that c eigenvalues are between 2 and 20 and all the others very close to zero (between 0.1 and 0.5) with c respectively equals to 5, 10 and 15. The residuals $\|Y - X\hat{\beta}_k\|^2$ are plotted for different values of k .

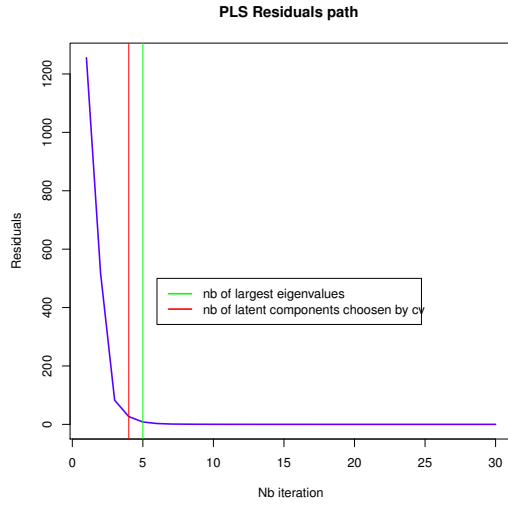


Figure 6

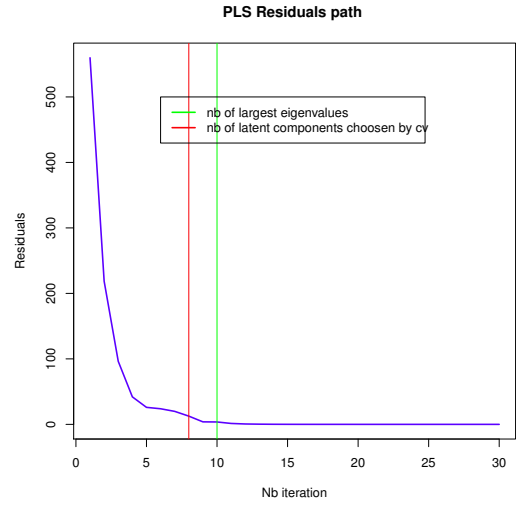


Figure 7

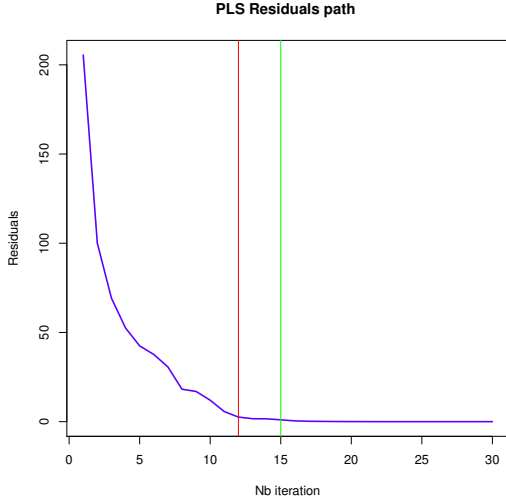


Figure 8

Figures above show that when there is a clear gap in the distribution of the eigenvalues with k large eigenvalues and the others very small then there is no need to go very far to recover most of the information with the PLS estimator. It is mainly due to the fact that the small eigenvalues can be considered as belonging to a same cluster.

5.2 Prediction error

Here we study the distance between the estimator and the true parameter in term of prediction error $\frac{1}{n} \|X\beta^* - X\hat{\beta}_k\|^2$. The expression of the prediction error is not as simple as the one of PCR and thus an upper bound for the prediction error is not as obvious since the PLS procedure is not a linear modeling procedure. Indeed, the direction of the new subspace onto which we project the observations depends in a complicated way on the singular value spectrum of the design matrix and also on the response. To compute or bound the prediction error we have to be careful because this also implies a control of the error due to the randomness of the subspace onto which we project your observations. We are going to use formula (5) of Theorem 4.1 to study the prediction error.

We first make the following assumptions. The real variables $\varepsilon_1, \dots, \varepsilon_n$ are assumed to be unobservable i.i.d centered gaussian random variables with common variance σ_n^2 . We also assume

- (H.1): $\sigma_n^2 = \mathcal{O}(\frac{1}{n})$. In other words we assume that $Y_i = x_i^T \beta^* + \delta_n \varepsilon_i$ where $\varepsilon_i \sim \mathcal{N}(0, 1)$ and $\delta_n = \frac{1}{\sqrt{n}}$ is the noise level which is related to the number of observations.
- (H.2): there exists a constant $L > 0$ such that $\min_{1 \leq i \leq n} \{p_i^2\} \geq L$, where $p_i = (X\beta^*)^T u_i$.

These two assumptions warrant that the signal to noise ratio $\left\{ \left| \frac{\tilde{\varepsilon}_i}{p_i} \right| \right\}_{1 \leq i \leq n}$ is not too small. This last quantity will appear again many time thereafter.

To bound from above the prediction error we have to be careful because PLS is a projected method but not onto a fixed subspace. The Krylov subspace onto which we project the data depends on Y and thus is a random subspace. Therefore we have to also control the randomness of the subspace onto which we project the data. To do so,

we introduce an oracle which is the regularization β_k of β^* onto the noise free Krylov subspace of dimension k . This regularized approximation of β^* is defined as

$$\beta_k \in \underset{\beta \in \mathcal{K}^k}{\operatorname{argmin}} \|X\beta^* - X\beta\|^2$$

where $\mathcal{K}^k := \mathcal{K}^k(X^T X, X^T X \beta^*)$ is the noise free Krylov subspace. Therefore we have

$$\beta_k = P_k^*(X^T X) X^T X \beta^*$$

where P_k^* is a polynomial of degree $k-1$ which satisfies

$$\|X\beta^* - X P_k^*(X^T X) X^T X \beta^*\|^2 = \underset{P \in \mathcal{P}_{k-1}}{\operatorname{argmin}} \|X\beta^* - X P(X^T X) X^T Y\|^2$$

and $X\beta^* - X\beta_k = Q_k^*(X X^T) X \beta^*$ with $Q_k^*(t) = 1 - t P_k^*(t) \in \mathcal{P}_{k,1}$. Then by the same arguments as the ones used to prove Proposition 3.2 and Theorem 4.1 we have that

1. the sequence of polynomials $(Q_k^*)_{1 \leq k \leq n}$ are orthogonals with respect to the measure

$$d\mu(\lambda) = \sum_{j=1}^n \lambda_j p_j^2 \delta_{\lambda_j},$$

where $p_j := u_j^T(X\beta^*)$.

- 2.

$$Q_k^*(\lambda_i) := \sum_{(j_1, \dots, j_k) \in I_k^+} w_{(j_1, \dots, j_k)} \prod_{l=1}^k \left(1 - \frac{\lambda_i}{\lambda_{j_l}}\right)$$

$$\text{where } w_{(j_1, \dots, j_k)} := \frac{p_{j_1}^2 \dots p_{j_k}^2 \lambda_{j_1}^2 \dots \lambda_{j_k}^2 V(\lambda_{j_1}, \dots, \lambda_{j_k})^2}{\sum_{(j_1, \dots, j_k) \in I_k^+} p_{j_1}^2 \dots p_{j_k}^2 \lambda_{j_1}^2 \dots \lambda_{j_k}^2 V(\lambda_{j_1}, \dots, \lambda_{j_k})^2}.$$

In fact we can write

$$\begin{aligned} & \frac{1}{n} \|X\beta^* - X\hat{\beta}_k\|^2 = \frac{1}{n} \|X\beta^* - X\hat{P}_k(X^T X) X^T Y\|^2 \\ & \leq \frac{2}{n} \|X\beta^* - X P_k^*(X^T X) X^T Y\|^2 + \frac{2}{n} \|X P_k^*(X^T X) X^T Y - X\hat{P}_k(X^T X) X^T Y\|^2 \\ & = \frac{2}{n} \|X\beta^* - X P_k^*(X^T X) X^T Y\|^2 + \frac{2}{n} \|\left(\hat{Q}_k(X X^T) - Q_k^*(X X^T)\right) Y\|^2. \end{aligned} \quad (9)$$

Therefore to bound by above $\frac{1}{n} \|X\beta^* - X\hat{\beta}_k\|^2$ we need to control two other quantities. The first one represents the error of regularization when projecting the linear predictor plus the noise on the observations onto the noise free Krylov subspace. The second quantities represents the approximation error between the projection onto the noise free Krylov subspace and onto the random Krylov subspace built from the observations.

Now let us introduce our main result on prediction error which provides an upper bound for the prediction error assuming a low variance of the observations noise.

Theorem 5.3. *Let $k < n$ and assume that (H.1) and (H.2) holds. Then, with probability at least $1 - n^{1-C}$ where $C > 1$, we have*

$$\frac{1}{n} \|X\beta^* - X\hat{\beta}_k\|^2 \leq$$

$$\begin{aligned} & \frac{1}{n} \left[2 \left(\frac{\sqrt{C(X^T X)} - 1}{\sqrt{C(X^T X)} + 1} \right)^{2k} + 4 \frac{\log(n)}{nL} \left(1 + \left(\frac{\sqrt{C(X^T X)} - 1}{\sqrt{C(X^T X)} + 1} \right)^{2k} \right) \right] \|X\beta^*\|^2 \\ & + \frac{4k^2 \tilde{C}^2 \log n}{L} \frac{1}{n^2} \left(1 + C \sqrt{\frac{\log n}{nL}} \right)^2 \|X\beta^*\|_W^2, \end{aligned}$$

where $C(X^T X) = \frac{\lambda_1}{\lambda_n}$, \tilde{C} is a constant and $W = \text{diag} \left(\max_{1 \leq i \leq n} \left(\prod_{l=1}^k \left| \frac{\lambda_i}{\lambda_{j_l}} - 1 \right|^2 \right) \right)$.

Proof. Theorem 5.3 is a straightforward consequence of Proposition 5.4, Proposition 5.5 and (9) below. In fact we recall that

$$\begin{aligned} & \frac{1}{n} \|X\beta^* - X\hat{\beta}_k\|^2 \\ & \leq \frac{2}{n} \|X\beta^* - XP_k^*(X^T X)X^T Y\|^2 + \frac{2}{n} \left\| \left(\hat{Q}_k(XX^T) - Q_k^*(XX^T) \right) Y \right\|^2. \end{aligned}$$

The following proposition provides an upper bound for the first term of (9).

Proposition 5.4. *With probability at least $1 - n^{1-C}$ where $C > 1$, we have for all $i = 1, \dots, n$*

$$\begin{aligned} & \frac{1}{n} \|X\beta^* - XP_k^*(X^T X)X^T Y\|^2 \\ & \leq \frac{1}{n} \left[2 \left(\frac{\sqrt{C(X^T X)} - 1}{\sqrt{C(X^T X)} + 1} \right)^{2k} + 4 \frac{\log(n)}{nL} \left(1 + \left(\frac{\sqrt{C(X^T X)} - 1}{\sqrt{C(X^T X)} + 1} \right)^{2k} \right) \right] \|X\beta^*\|^2. \end{aligned}$$

Then we bound by above the second term $\frac{1}{n} \left\| \left(\hat{Q}_k(XX^T) - Q_k^*(XX^T) \right) Y \right\|^2$.

Proposition 5.5. *Assume (H.1) and (H.2). Then with probability larger than $1 - n^{1-C}$ where $C > 1$ we have*

$$\begin{aligned} & \frac{1}{n} \left\| \left(\hat{Q}_k(X^T X) - Q_k^*(X^T X) \right) Y \right\|^2 \\ & \leq \frac{4k^2 \tilde{C}^2 \log n}{L} \frac{1}{n^2} \left(1 + C \sqrt{\frac{\log n}{nL}} \right)^2 \|X\beta^*\|_W^2, \end{aligned}$$

where $W = \text{diag} \left(\max_{1 \leq i \leq n} \left(\prod_{l=1}^k \left| \frac{\lambda_i}{\lambda_{j_l}} - 1 \right|^2 \right) \right)$ and \tilde{C} is a constant.

The theorem is proved by combining the two previous bounds. □

The bound in Theorem 5.3 highly depends on the signal to noise ratio which must not be too small with respect to the eigenvector directions of $X^T X$ to ensure good statistical properties of the PLS estimator. This is the major difference with PCA that takes into account the variance of the noise on the observations to build the latent variables but not the level of the signal. On the contrary PLS takes into account the signal through Y to construct the latent variables. That is why for PLS the signal to noise ratio plays an important role in the accuracy of the model.

The following simulation highlights this statement showing that there is generally no hope to recover a good approximation of the predicted function in case of a high variance

of the noise. We have compared the performances of the PLS estimator at different steps and for different levels of noise. Here is a typical example of the behaviour of the prediction error for the PLS estimator. The data set is simulated according to model (1) with $n = p = 100$. Figure 9 represents the PLS prediction error path for different value of the parameter k and for different level of noise on the observations.

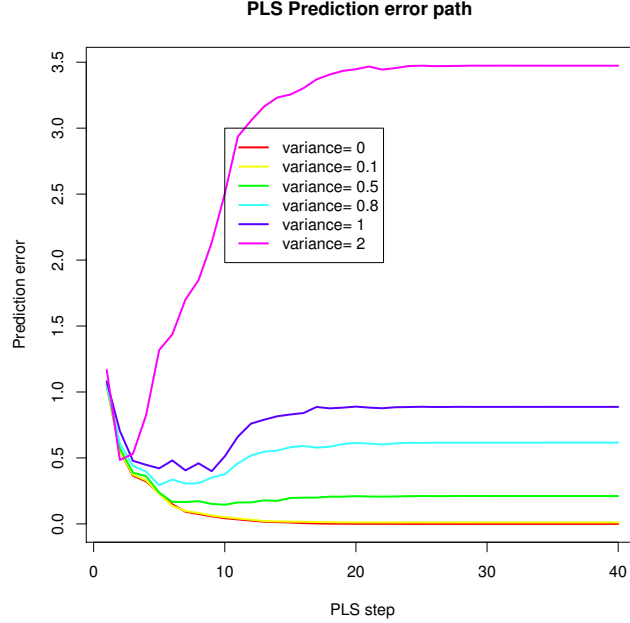


Figure 9

Figures above show that there is no assurance that the prediction error goes to zero when the variance is too high compared to the number of observations. This is essentially due to the fact that the PLS method is an iterative technique and thus the noise can propagate at each step of the construction leading to an undue amplification of the error.

6 Conclusion

PLS is a method used to remove multicollinearities and based on the construction of a new subspace of reduced dimension. This new subspace is built to maximize both the covariance of the covariates and the correlation to the response. The key idea behind PLS is to approximate and regularize the pseudo-inverse of the covariance matrix by a polynomial in the power of the matrix. PLS is in fact a least square problem with random linear constraints. This method can also be viewed as a minimization problem over a particular polynomial subspace. From this perspective we showed that the PLS residuals are in fact orthogonal polynomials with respect to a measure based on the spectrum of the covariance matrix. From the definition of this discrete measure we deduce a new formula for the residuals. This formula depends explicitly on the observations noise and on the spectrum of the covariance matrix. At last we have taken advantages of these findings in a regression context to state new results for the estimation and prediction error for PLS under a low variance of the noise. The control of the signal-to-noise ratio and of the spectrum distribution seems to be the key to state such results. We have showed that PLS is not an automatic solution to avoid the problem of multicollinearity in regression. We

have to be careful when using PLS because its statistical properties are strongly depending on the features of the data and in particular on the distribution of the spectrum. The main drawback of PLS is the fact that it seems inappropriate if Y has too much variation but the advantages is that it takes both X and Y into account in the decomposition of X contrary to PCR. To conclude this paper throw new lights on PLS in a regression context but this is not the end of the road and the formula for the residuals should be explored further to completely understand the method.

7 Proof

7.1 Proof of Proposition 3.2

Proof. Let $k \in \mathbb{N}^*$ and $n > l > k$. Because $\hat{Q}_k \in \mathcal{P}_{k,1}$ we have

$$XX^T \hat{Q}_k(XX^T)Y \in \mathcal{K}^{k+1}(XX^T, XX^T Y).$$

Furthermore from (4) we get $\hat{Q}_l(XX^T)Y \perp \mathcal{K}^l(XX^T, XX^T Y)$. Besides $\mathcal{K}^l(XX^T, XX^T Y) \supset \mathcal{K}^{k+1}(XX^T, XX^T Y)$. Therefore we deduce that for all $k \neq l$ we have $X^T \hat{Q}_k(XX^T)Y \perp \hat{Q}_l(XX^T)Y$. Then using the SVD decomposition of X we get $XX^T = \sum_{1 \leq j \leq n} \lambda_j u_j u_j^T$ and

$$\begin{aligned} 0 &= \left\langle XX^T \hat{Q}_k(XX^T)Y, \hat{Q}_l(XX^T)Y \right\rangle \\ &= \left(\sum_{1 \leq j \leq n} \lambda_j \hat{Q}_k(\lambda_j) u_j u_j^T Y \right)^T \left(\sum_{1 \leq j \leq n} \hat{Q}_l(\lambda_j) u_j u_j^T Y \right) = \sum_{1 \leq j \leq n} \lambda_j \hat{Q}_k(\lambda_j) \hat{Q}_l(\lambda_j) (u_j^T Y)^2. \end{aligned}$$

Finally we get

$$0 = \sum_{1 \leq j \leq n} \lambda_j \hat{Q}_k(\lambda_j) \hat{Q}_l(\lambda_j) (u_j^T Y)^2.$$

And we deduce that $\hat{Q}_1, \hat{Q}_2, \dots, \hat{Q}_{n-1}$ is a sequence of orthonormal polynomials with respect to the measure

$$d\hat{\mu}(\lambda) = \sum_{j=1}^n \lambda_j (u_j^T Y)^2 \delta_{\lambda_j}.$$

□

7.2 Proof of Theorem 4.1

We recall that $(\hat{Q}_1)_{1 \leq k < n}$ is a sequence of orthonormal polynomials with respect to the measure $d\hat{\mu}(\lambda)$. Returning to the definition of orthogonal polynomials we first express the polynomials $(\hat{Q}_k)_{1 \leq k < n}$ as the quotient of two determinants.

Proposition 7.1. *For all $j \in \mathbb{N}$, let $\hat{m}_j = \int x^j d\hat{\mu}$.*

Then for all $k \in \llbracket 1, \dots, n-1 \rrbracket$ we have

$$\hat{Q}_k(x) = (-1)^k \frac{\det(\hat{G}_{2k-1}(x))}{\det(\hat{H}_{2k-1})} \quad (10)$$

where

$$\hat{G}_{2k-1}(x) := \begin{bmatrix} \hat{m}_0 & \hat{m}_1 & \dots & \hat{m}_k \\ \vdots & & & \\ \hat{m}_{k-1} & \hat{m}_k & \dots & \hat{m}_{2k-1} \\ 1 & x & \dots & x^k \end{bmatrix}$$

and

$$\hat{H}_{2k-1} := \begin{bmatrix} \hat{m}_1 & \hat{m}_2 & \dots & \hat{m}_k \\ \vdots & & & \\ \hat{m}_{k-1} & \hat{m}_k & \dots & \hat{m}_{2k-2} \\ \hat{m}_k & \hat{m}_{k+1} & \dots & \hat{m}_{2k-1} \end{bmatrix}.$$

Proof. The polynomials $(\hat{Q}_k)_{1 \leq k < n}$ are the ones which satisfy

1. $\hat{Q}_k(x) = \alpha_k^k x^k + \alpha_k^{k-1} x^{k-1} + \dots + \alpha_1^k x + \alpha_0^k$
2. $\forall j \in [0, k-1], \int \left[x^j (\alpha_k^k x^k + \alpha_k^{k-1} x^{k-1} + \dots + \alpha_1^k x + \alpha_0^k) \right] d\hat{\mu} = 0$
3. $\hat{Q}_k(0) = 1$

This is equivalent to solve the following system of k equations with k unknowns

$$\forall j \in \llbracket 0, k-1 \rrbracket, \quad \alpha_k^k \hat{m}_{j+k} + \alpha_k^{k-1} \hat{m}_{j+k-1} + \dots + \alpha_1^k \hat{m}_{j+1} = -\hat{m}_j.$$

The solution $(\alpha_1^k, \dots, \alpha_k^k)$ of this system satisfies

$$\begin{bmatrix} \hat{m}_1 & \hat{m}_2 & \dots & \hat{m}_k \\ \vdots & & & \\ \hat{m}_{k-1} & \hat{m}_k & \dots & \hat{m}_{2k-2} \\ \hat{m}_k & \hat{m}_{k+1} & \dots & \hat{m}_{2k-1} \end{bmatrix} \begin{bmatrix} \alpha_1^k \\ \alpha_2^k \\ \vdots \\ \alpha_k^k \end{bmatrix} = - \begin{bmatrix} \hat{m}_0 \\ \hat{m}_1 \\ \vdots \\ \hat{m}_{k-1} \end{bmatrix}$$

We conclude the proof using the Cramer's rule which provides explicit formula for the solution of a system of linear equations with as many equations as unknowns. \square

Then returning to the definition of the discrete measure $\hat{\mu}$ we explicitly express $\hat{Q}_k(\lambda_i)$ in terms of $(\lambda_i)_{1 \leq i \leq n}$ and $(u_i^T Y)_{1 \leq i \leq n}$ for all $1 \leq k < n$ and all $1 \leq i \leq n$.

Proposition 7.2. *Let $k \in \llbracket 1, \dots, n-1 \rrbracket$ and $i \in \llbracket 1, \dots, n \rrbracket$.*

Let $\hat{p}_i := Y^T u_i$. Define

$$I_k = \left\{ (j_1, \dots, j_k) \in \llbracket 1, n \rrbracket^k, j_1 \neq \dots \neq j_k \right\}$$

and

$$I_{k,i} = \left\{ (j_1, \dots, j_k) \in \llbracket 1, n \rrbracket^k, j_1 \neq \dots \neq j_k \neq i \right\}.$$

We have

$$\hat{Q}_k(\lambda_i) = (-1)^k \frac{\sum_{(j_1, \dots, j_k) \in I_{k,i}} \hat{p}_{j_1}^2 \dots \hat{p}_{j_k}^2 V(\lambda_{j_1}, \dots, \lambda_{j_k}, \lambda_i) \lambda_{j_1} \dots \lambda_{j_k}^k}{\sum_{(j_1, \dots, j_k) \in I_k} \hat{p}_{j_1}^2 \dots \hat{p}_{j_k}^2 V(\lambda_{j_1}, \dots, \lambda_{j_k}) \lambda_{j_1}^2 \dots \lambda_{j_k}^{k+1}} \quad (11)$$

where $V(x_1, \dots, x_l)$ is the Vandermonde determinant of $(x_1, \dots, x_l) \in \mathbb{R}^l$.

If $k = n$ we have

$$\hat{Q}_k(\lambda_i) = 0.$$

Proof. Let $1 \leq i \leq n$. Using the fact that $d\hat{\mu} = \sum_{j=1}^n \lambda_j \hat{p}_j^2 \delta_{\lambda_j}$ we get

$$\begin{aligned}
& \det \begin{bmatrix} \hat{m}_0 & \hat{m}_1 & \dots & \hat{m}_k \\ \vdots & & & \\ \hat{m}_{k-1} & \hat{m}_k & & \hat{m}_{2k-1} \\ 1 & \lambda_i & \dots & \lambda_i^k \end{bmatrix} \\
&= \det \begin{bmatrix} \sum_{j=1}^n \lambda_j \hat{p}_j^2 & \sum_{j=1}^n \lambda_j^2 \hat{p}_j^2 & \dots & \sum_{j=1}^n \lambda_j^{k+1} \hat{p}_j^2 \\ \vdots & & & \\ \sum_{j=1}^n \lambda_j^k \hat{p}_j^2 & \sum_{j=1}^n \lambda_j^2 \hat{p}_j^2 & & \sum_{j=1}^n \lambda_j^{2k} \hat{p}_j^2 \\ 1 & \lambda_i & \dots & \lambda_i^k \end{bmatrix} \\
&= \sum_{j_1=1}^n \dots \sum_{j_k=1}^n \hat{p}_{j_1}^2 \hat{p}_{j_2}^2 \dots \hat{p}_{j_k}^2 \lambda_{j_1} \lambda_{j_2}^2 \dots \lambda_{j_k}^k \det \begin{bmatrix} 1 & \lambda_{j_1} & \dots & \lambda_{j_1}^k \\ \vdots & & & \\ 1 & \lambda_{j_k} & & \lambda_{j_k}^k \\ 1 & \lambda_i & \dots & \lambda_i^k \end{bmatrix}
\end{aligned}$$

where

$$\det \begin{bmatrix} 1 & \lambda_{j_1} & & \lambda_{j_1}^k \\ \vdots & & & \\ 1 & \lambda_{j_k} & & \lambda_{j_k}^k \\ 1 & \lambda_i & \dots & \lambda_i^k \end{bmatrix} = V(\lambda_{j_1}, \dots, \lambda_{j_k}, \lambda_i).$$

$V(\lambda_{j_1}, \dots, \lambda_{j_k}, \lambda_i)$ is the Vandermonde determinant of $\lambda_{j_1}, \dots, \lambda_{j_k}, \lambda_i$ and is non zero only if all the $\lambda_{j_1}, \dots, \lambda_{j_k}, \lambda_i$ are distincts.

Therefore if $k < n$, we get

$$\det \begin{bmatrix} \hat{m}_0 & \hat{m}_1 & \dots & \hat{m}_k \\ \vdots & & & \\ \hat{m}_{k-1} & \hat{m}_k & & \hat{m}_{2k-1} \\ 1 & \lambda_i & \dots & \lambda_i^k \end{bmatrix} = \sum_{(j_1, \dots, j_k) \in I_{k,i}} \hat{p}_{j_1}^2 \hat{p}_{j_2}^2 \dots \hat{p}_{j_k}^2 \lambda_{j_1} \lambda_{j_2}^2 \dots \lambda_{j_k}^k V(\lambda_{j_1}, \dots, \lambda_{j_k}, \lambda_i). \quad (12)$$

Using the same arguments we also get

$$\det \begin{bmatrix} \hat{m}_1 & \hat{m}_2 & \dots & \hat{m}_k \\ \vdots & & & \\ \hat{m}_{k-1} & \hat{m}_k & & \hat{m}_{2k-2} \\ \hat{m}_k & \hat{m}_{k+1} & \dots & \hat{m}_{2k-1} \end{bmatrix} = \sum_{(j_1, \dots, j_k) \in I_k} \hat{p}_{j_1}^2 \hat{p}_{j_2}^2 \dots \hat{p}_{j_k}^2 \lambda_{j_1}^2 \lambda_{j_2}^3 \dots \lambda_{j_k}^{k+1} V(\lambda_{j_1}, \dots, \lambda_{j_k}). \quad (13)$$

From (10), (12) and (13) we deduce (11).

$$\begin{aligned}
& \text{When } k = n, \det \begin{bmatrix} \hat{m}_0 & \hat{m}_1 & \dots & \hat{m}_k \\ \vdots & & & \\ \hat{m}_{k-1} & \hat{m}_k & & \hat{m}_{2k-1} \\ 1 & \lambda_i & \dots & \lambda_i^k \end{bmatrix} = 0 \text{ and therefore} \\
& \hat{Q}_k(\lambda_i) = 0.
\end{aligned}$$

□

Now using the properties of the Vandermonde determinant we provide a more useful characterization of the residual $\hat{Q}_k(\lambda_i)$. Let $k < n$. Formula (11) of Proposition 7.2 tells us that

$$\hat{Q}_k(\lambda_i) = (-1)^k \frac{\sum_{(j_1, \dots, j_k) \in I_{k,i}} \hat{p}_{j_1}^2 \dots \hat{p}_{j_k}^2 V(\lambda_{j_1}, \dots, \lambda_{j_k}, \lambda_i) \lambda_{j_1} \dots \lambda_{j_k}^k}{\sum_{(j_1, \dots, j_k) \in I_k} \hat{p}_{j_1}^2 \dots \hat{p}_{j_k}^2 V(\lambda_{j_1}, \dots, \lambda_{j_k}) \lambda_{j_1}^2 \dots \lambda_{j_k}^{k+1}}. \quad (14)$$

On the one hand, we have

$$\begin{aligned} & \sum_{(j_1, \dots, j_k) \in I_k} \hat{p}_{j_1}^2 \dots \hat{p}_{j_k}^2 V(\lambda_{j_1}, \dots, \lambda_{j_k}) \lambda_{j_1}^2 \dots \lambda_{j_k}^{k+1} \\ &= \sum_{(j_1, \dots, j_k) \in I_k^+} \sum_{\tau \in \mathcal{S}(1, \dots, k)} \hat{p}_{j_{\tau(1)}}^2 \dots \hat{p}_{j_{\tau(k)}}^2 V(\lambda_{j_{\tau(1)}}, \dots, \lambda_{j_{\tau(k)}}) \lambda_{j_{\tau(1)}}^2 \dots \lambda_{j_{\tau(k)}}^{k+1} \end{aligned}$$

where $\mathcal{S}(1, \dots, k)$ is the set formed of all the permutations of $(1, \dots, k)$. Then using the fact that $V(\lambda_{j_{\tau(1)}}, \dots, \lambda_{j_{\tau(k)}}) = \varepsilon(\tau) V(\lambda_{j_1}, \dots, \lambda_{j_k})$ we get

$$\begin{aligned} & \sum_{(j_1, \dots, j_k) \in I_k} \hat{p}_{j_1}^2 \dots \hat{p}_{j_k}^2 V(\lambda_{j_1}, \dots, \lambda_{j_k}) \lambda_{j_1}^2 \dots \lambda_{j_k}^{k+1} \\ &= \sum_{(j_1, \dots, j_k) \in I_k^+} \sum_{\tau \in \mathcal{S}(1, \dots, k)} \hat{p}_{j_1}^2 \dots \hat{p}_{j_k}^2 \varepsilon(\tau) V(\lambda_{j_1}, \dots, \lambda_{j_k}) \lambda_{j_1}^2 \dots \lambda_{j_k}^2 \lambda_{j_{\tau(2)}} \dots \lambda_{j_{\tau(k)}}^{k-1} \\ &= \sum_{(j_1, \dots, j_k) \in I_k^+} \hat{p}_{j_1}^2 \dots \hat{p}_{j_k}^2 V(\lambda_{j_1}, \dots, \lambda_{j_k}) \lambda_{j_1}^2 \dots \lambda_{j_k}^2 \left[\sum_{\tau \in \mathcal{S}(1, \dots, k)} \varepsilon(\tau) \lambda_{j_{\tau(2)}} \dots \lambda_{j_{\tau(k)}}^{k-1} \right]. \quad (15) \end{aligned}$$

On the other hand,

$$V(\lambda_{j_1}, \dots, \lambda_{j_k}) = \sum_{\tau \in \mathcal{S}(1, \dots, k)} \varepsilon(\tau) \lambda_{j_1}^{\tau(1)-1} \dots \lambda_{j_k}^{\tau(k)-1} = \sum_{\tau \in \mathcal{S}(1, \dots, k)} \varepsilon(\tau) \lambda_{j_{\tau(2)}} \dots \lambda_{j_{\tau(k)}}^{k-1}. \quad (16)$$

To conclude (15) and (16) leads to

$$\sum_{(j_1, \dots, j_k) \in I_k} \hat{p}_{j_1}^2 \dots \hat{p}_{j_k}^2 V(\lambda_{j_1}, \dots, \lambda_{j_k}) \lambda_{j_1}^2 \dots \lambda_{j_k}^{k+1} = \sum_{(j_1, \dots, j_k) \in I_k^+} \hat{p}_{j_1}^2 \dots \hat{p}_{j_k}^2 \lambda_{j_1}^2 \dots \lambda_{j_k}^2 V(\lambda_{j_1}, \dots, \lambda_{j_k})^2. \quad (17)$$

A similar reasoning can be applied to the numerator. Indeed using the fact that

$$V(\lambda_{j_1}, \dots, \lambda_{j_k}, \lambda_i) = \prod_{l=1}^k (\lambda_i - \lambda_{j_l}) \prod_{1 \leq q < m \leq k} (\lambda_{j_m} - \lambda_{j_q}) = \prod_{l=1}^k (\lambda_i - \lambda_{j_l}) V(\lambda_{j_1}, \dots, \lambda_{j_k})$$

we get

$$\begin{aligned} & \sum_{(j_1, \dots, j_k) \in I_{k,i}} \hat{p}_{j_1}^2 \dots \hat{p}_{j_k}^2 V(\lambda_{j_1}, \dots, \lambda_{j_k}, \lambda_i) \lambda_{j_1} \dots \lambda_{j_k}^k \\ &= \sum_{(j_1, \dots, j_k) \in I_{k,i}^+} \sum_{\tau \in \mathcal{S}(1, \dots, k)} \hat{p}_{j_{\tau(1)}}^2 \dots \hat{p}_{j_{\tau(k)}}^2 \prod_{l=1}^k (\lambda_i - \lambda_{j_{\tau(l)}}) V(\lambda_{j_{\tau(1)}}, \dots, \lambda_{j_{\tau(k)}}) \lambda_{j_{\tau(1)}} \dots \lambda_{j_{\tau(k)}}^k \end{aligned}$$

$$\begin{aligned}
&= \sum_{(j_1, \dots, j_k) \in I_{k,i}^+} \hat{p}_{j_1}^2 \dots \hat{p}_{j_k}^2 \prod_{l=1}^k (\lambda_i - \lambda_{j_l}) V(\lambda_{j_1}, \dots, \lambda_{j_k}) \lambda_{j_1} \dots \lambda_{j_k} \left[\sum_{\tau \in \mathcal{S}(1, \dots, k)} \varepsilon(\tau) \lambda_{j_{\tau(2)}} \dots \lambda_{j_{\tau(k)}}^{k-1} \right] \\
&= \sum_{(j_1, \dots, j_k) \in I_k^+} \hat{p}_{j_1}^2 \dots \hat{p}_{j_k}^2 \lambda_{j_1} \dots \lambda_{j_k} V(\lambda_{j_1}, \dots, \lambda_{j_k})^2 \prod_{l=1}^k (\lambda_i - \lambda_{j_l}) \\
&= (-1)^k \sum_{(j_1, \dots, j_k) \in I_k^+} \hat{p}_{j_k}^2 \lambda_{j_1}^2 \dots \lambda_{j_k}^2 V(\lambda_{j_1}, \dots, \lambda_{j_k})^2 \prod_{l=1}^k (1 - \frac{\lambda_i}{\lambda_{j_l}}). \tag{18}
\end{aligned}$$

From (14), (17) and (18) we conclude

$$\hat{Q}_k(\lambda_i) = \sum_{(j_1, \dots, j_k) \in I_k^+} \left[\frac{\hat{p}_{j_1}^2 \dots \hat{p}_{j_k}^2 \lambda_{j_1}^2 \dots \lambda_{j_k}^2 V(\lambda_{j_1}, \dots, \lambda_{j_k})^2}{\sum_{(j_1, \dots, j_k) \in I_k^+} \hat{p}_{j_1}^2 \dots \hat{p}_{j_k}^2 \lambda_{j_1}^2 \dots \lambda_{j_k}^2 V(\lambda_{j_1}, \dots, \lambda_{j_k})^2} \right] \prod_{l=1}^k (1 - \frac{\lambda_i}{\lambda_{j_l}}).$$

7.3 Proof of Proposition 5.1

Proof. Let $k < n$. By definition of $\hat{\beta}_k$ and referring to results of Proposition 3.1 we have

$$\|Y - X\hat{\beta}_k\|^2 = \min_{Q \in \mathcal{P}_{k,1}} \|Q(XX^T)Y\|^2 = \min_{Q \in \mathcal{P}_{k,1}} \|Q(XX^T)(X\beta^* + \varepsilon)\|^2.$$

Using the decomposition of β^* and ε on the left and right eigenvectors (i.e. $\beta^* = \sum_{i=1}^p \tilde{\beta}_i^* v_i$ where $\tilde{\beta}_i^* = \beta^T v_i$ and $\varepsilon = \sum_{i=1}^n \tilde{\varepsilon}_i u_i$ where $\tilde{\varepsilon}_i = \varepsilon^T u_i$) we get

$$\begin{aligned}
\|Y - X\hat{\beta}_k\|^2 &= \min_{Q \in \mathcal{P}_{k,1}} \left(\sum_{i=1}^n Q(\lambda_i)^2 \left(\sqrt{\lambda_i} \tilde{\beta}_i^* + \tilde{\varepsilon}_i \right)^2 \right) \\
&\leq \left(\min_{Q \in \mathcal{P}_{k,1}} \max_{\lambda \in [\lambda_n, \lambda_1]} Q(\lambda)^2 \right) \sum_{i=1}^n \left(\sqrt{\lambda_i} \tilde{\beta}_i^* + \tilde{\varepsilon}_i \right)^2.
\end{aligned}$$

Then we have

$$\begin{aligned}
\|Y - X\hat{\beta}_k\|^2 &\leq \left(\min_{Q \in \mathcal{P}_{k,1}} \max_{\lambda \in [\lambda_n, \lambda_1]} |Q(\lambda)| \right)^2 \sum_{i=1}^n \left(\sqrt{\lambda_i} \tilde{\beta}_i^* + \tilde{\varepsilon}_i \right)^2 \\
&\leq \frac{1}{\left[C_k \left(\frac{\lambda_1 + \lambda_n}{\lambda_1 - \lambda_n} \right) \right]^2} \sum_{i=1}^n \left(\sqrt{\lambda_i} \tilde{\beta}_i^* + \tilde{\varepsilon}_i \right)^2
\end{aligned}$$

where C_k is the k^{th} Chebyshev polynomial. This last inequalities follows from Proposition 5.2. Then we use the fact that

$$\begin{aligned}
\left| C_k \left(\frac{\lambda_1 + \lambda_n}{\lambda_1 - \lambda_n} \right) \right| &= \frac{1}{2} \left| \left(\frac{\sqrt{\lambda_1} + \sqrt{\lambda_n}}{\sqrt{\lambda_1} - \sqrt{\lambda_n}} \right)^k + \left(\frac{\sqrt{\lambda_1} - \sqrt{\lambda_n}}{\sqrt{\lambda_1} + \sqrt{\lambda_n}} \right)^k \right| \\
&= \frac{1}{2} \left| \left(\frac{\sqrt{C(X^T X)} + 1}{\sqrt{C(X^T X)} - 1} \right)^k + \left(\frac{\sqrt{C(X^T X)} - 1}{\sqrt{C(X^T X)} + 1} \right)^k \right| \geq \left(\frac{\sqrt{C(X^T X)} + 1}{\sqrt{C(X^T X)} - 1} \right)^k,
\end{aligned}$$

where $C(X^T X) = \frac{\lambda_1}{\lambda_n}$. At last we get

$$\mathbb{E} \left[\frac{1}{n} \|Y - X\hat{\beta}_k\|^2 \right] \leq \frac{1}{n} \left(\frac{\sqrt{C(X^T X)} - 1}{\sqrt{C(X^T X)} + 1} \right)^{2k} \mathbb{E} \left[\sum_{i=1}^n \left(\sqrt{\lambda_i} \tilde{\beta}_i^* + \tilde{\varepsilon}_i \right)^2 \right]$$

and since the $(\varepsilon_j)_{1 \leq j \leq n}$ are assumed to be centered we conclude

$$\mathbb{E} \left[\frac{1}{n} \|Y - X\hat{\beta}_k\|^2 \right] \leq \left(\frac{\sqrt{C(X^T X)} - 1}{\sqrt{C(X^T X)} + 1} \right)^{2k} \left[\frac{1}{n} \|X\beta^*\|^2 + \sigma^2 \right].$$

□

7.4 Proof of Proposition 5.4

Proof. We have

$$\begin{aligned} & \frac{1}{n} \|X\beta^* - XP_k^*(X^T X)X^T Y\|^2 \\ & \leq \frac{2}{n} \|X\beta^* - XP_k^*(X^T X)X^T X^T X\beta^*\|^2 + \frac{2}{n} \|XP_k^*(X^T X)X^T \varepsilon\|^2 \\ & = \frac{2}{n} \|Q_k^*(XX^T)X\beta^*\|^2 + \frac{2}{n} \|XP_k^*(X^T X)X^T \varepsilon\|^2. \end{aligned}$$

On one hand, by the same arguments as the ones used to prove Proposition 5.1 (with no noise), we get

$$\frac{1}{n} \|Q_k^*(XX^T)X\beta^*\|^2 \leq \frac{1}{n} \left(\frac{\sqrt{C(X^T X)} - 1}{\sqrt{C(X^T X)} + 1} \right)^k \|X\beta^*\|^2 \quad (19)$$

where $C(X^T X) = \frac{\lambda_1}{\lambda_n}$ is the ratio of the two extreme non zero eigenvalues of $X^T X$.

On the other hand we have

$$\begin{aligned} \frac{1}{n} \|XP_k^*(X^T X)X^T \varepsilon\|^2 &= \frac{1}{n} \sum_{i=1}^n (1 - Q_k^*(\lambda_i))^2 \tilde{\varepsilon}_i^2 \\ &= \frac{1}{n} \sum_{i=1}^n (1 - Q_k^*(\lambda_i))^2 p_i^2 \frac{\tilde{\varepsilon}_i^2}{p_i^2} \end{aligned}$$

where $\tilde{\varepsilon}_i = \varepsilon^T u_i$. Notice that $Q_k^*(\lambda_i)$ can be positive or negative and therefore the factors in the last sum oscillate above and below one (see Subsection 4.2). We bound this last term by above using concentration inequalities. Here a low variance of the noise is necessary to ensure that the term we consider is not too large. The random variables $(\varepsilon_i)_{1 \leq i \leq n}$ are assumed to be i.i.d $\sim \mathcal{N}(0, \sigma_n^2)$ and so are the $(\tilde{\varepsilon}_i)_{1 \leq i \leq n}$. Therefore we use the following proposition which is a direct consequence of concentration inequalities for Gaussian random variables

Proposition 7.3. *Let $\mathcal{A} = \{\cap_{i=1}^n |\tilde{\varepsilon}_i| \leq \delta\}$. If assumptions (H.1) holds then there exists a constant $C > 1$ such that*

$$\mathbb{P}(\mathcal{A}^c) \leq \sum_{i=1}^n \mathbb{P}(|\tilde{\varepsilon}_i| > \delta) \leq \sum_{i=1}^n e^{-\frac{\delta^2}{2\sigma_n^2}} \leq ne^{-C\delta^2 n}.$$

In addition with probability at least $1 - n^{1-C}$ we have for all $i = 1, \dots, n$

$$|\tilde{\varepsilon}_i| \leq \sqrt{\frac{\log(n)}{n}}$$

With Proposition 7.3 we deduce that with probability at least $1 - n^{1-C}$ where $C > 1$ we have for all $i = 1, \dots, n$

$$\frac{1}{n} \|XP_k^*(X^T X)X^T \varepsilon\|^2 \leq \frac{1}{n} \left(\sum_{i=1}^n (1 - Q_k^*(\lambda_i))^2 p_i^2 \right) \frac{\log(n)}{nL}.$$

Then using the triangular inequality and (19) we state

$$\frac{1}{n} \|XP_k^*(X^T X)X^T \varepsilon\|^2 \leq \frac{1}{n} \left[2 + 2 \left(\frac{\sqrt{C(X^T X)} - 1}{\sqrt{C(X^T X)} + 1} \right)^{2k} \right] \|X\beta^*\|^2 \frac{\log(n)}{nL}. \quad (20)$$

Combining (19) and (20) we conclude

$$\begin{aligned} & \frac{1}{n} \|X\beta^* - XP_k^*(X^T X)X^T Y\|^2 \\ & \leq \frac{1}{n} \left[2 \left(\frac{\sqrt{C(X^T X)} - 1}{\sqrt{C(X^T X)} + 1} \right)^{2k} + 4 \frac{\log(n)}{nL} \left(1 + \left(\frac{\sqrt{C(X^T X)} - 1}{\sqrt{C(X^T X)} + 1} \right)^{2k} \right) \right] \|X\beta^*\|^2. \end{aligned} \quad (21)$$

□

7.5 Proof of Proposition 5.5

Proof. Using the SVD of XX^T we get

$$\frac{1}{n} \left\| \left(\hat{Q}_k(X^T X) - Q_k^*(X^T X) \right) Y \right\|^2 = \frac{1}{n} \sum_{i=1}^n \left(\hat{Q}_k(\lambda_i) - Q_k^*(\lambda_i) \right)^2 \hat{p}_i^2. \quad (22)$$

Remark We can notice that

$$\sum_{i=1}^n \left(\hat{Q}_k(\lambda_i) - Q_k^*(\lambda_i) \right)^2 \hat{p}_i^2 \leq \frac{1}{\lambda_n} \sum_{i=1}^n \left(\hat{Q}_k(\lambda_i) - Q_k^*(\lambda_i) \right)^2 \lambda_i \hat{p}_i^2 \leq \frac{1}{\lambda_n} \|\hat{Q}_k - Q_k^*\|_{\hat{\mu}}^2.$$

We define

$$\hat{D}_{j_1, \dots, j_k} := \hat{p}_{j_1}^2 \dots \hat{p}_{j_k}^2 \lambda_{j_1}^2 \dots \lambda_{j_k}^2 V(\lambda_{j_1}, \dots, \lambda_{j_k})^2 > 0,$$

$$D_{j_1, \dots, j_k} := p_{j_1}^2 \dots p_{j_k}^2 \lambda_{j_1}^2 \dots \lambda_{j_k}^2 V(\lambda_{j_1}, \dots, \lambda_{j_k})^2 > 0.$$

and

$$\hat{D}_k := \sum_{(j_1, \dots, j_k) \in I_k^+} \hat{D}_{j_1, \dots, j_k}$$

,

$$D_k := \sum_{(j_1, \dots, j_k) \in I_k^+} D_{j_1, \dots, j_k}.$$

We recall that

$$\hat{Q}_k(\lambda_i) = (-1)^k \frac{\sum_{(j_1, \dots, j_k) \in I_k^+} \hat{D}_{j_1, \dots, j_k} \prod_{l=1}^k \left(\frac{\lambda_i}{\lambda_{j_l}} - 1 \right)}{\sum_{(j_1, \dots, j_k) \in I_k^+} \hat{D}_{j_1, \dots, j_k}}$$

and

$$Q_k(\lambda_i) = (-1)^k \frac{\sum_{(j_1, \dots, j_k) \in I_k^+} D_{j_1, \dots, j_k} \prod_{l=1}^k \left(\frac{\lambda_i}{\lambda_{j_l}} - 1 \right)}{\sum_{(j_1, \dots, j_k) \in I_k^+} D_{j_1, \dots, j_k}}.$$

We have

$$\begin{aligned}
\left| \hat{Q}_k(\lambda_i) - Q_k^*(\lambda_i) \right| &\leq \left| \frac{\sum_{(j_1, \dots, j_k) \in I_k^+} \left[\hat{D}_{j_1, \dots, j_k} \prod_{l=1}^k \left(\frac{\lambda_i}{\lambda_{j_l}} - 1 \right) \right]}{\hat{D}_k} - \frac{\sum_{(j_1, \dots, j_k) \in I_k^+} \left[D_{j_1, \dots, j_k} \prod_{l=1}^k \left(\frac{\lambda_i}{\lambda_{j_l}} - 1 \right) \right]}{D_k} \right| \\
&\leq \left| \frac{\sum_{(j_1, \dots, j_k) \in I_k^+} \left[D_{j_1, \dots, j_k} \prod_{l=1}^k \left(\frac{\lambda_i}{\lambda_{j_l}} - 1 \right) \right]}{D_k} - \frac{\sum_{(j_1, \dots, j_k) \in I_k^+} \left[\hat{D}_{j_1, \dots, j_k} \prod_{l=1}^k \left(\frac{\lambda_i}{\lambda_{j_l}} - 1 \right) \right]}{D_k} \right| \\
&\quad + \left| \frac{\sum_{(j_1, \dots, j_k) \in I_k^+} \left[\hat{D}_{j_1, \dots, j_k} \prod_{l=1}^k \left(\frac{\lambda_i}{\lambda_{j_l}} - 1 \right) \right]}{D_k} - \frac{\sum_{(j_1, \dots, j_k) \in I_k^+} \left[\hat{D}_{j_1, \dots, j_k} \prod_{l=1}^k \left(\frac{\lambda_i}{\lambda_{j_l}} - 1 \right) \right]}{\hat{D}_k} \right| \\
&\leq \frac{1}{D_k} \left| \sum_{(j_1, \dots, j_k) \in I_k^+} \left[D_{j_1, \dots, j_k} - \hat{D}_{j_1, \dots, j_k} \right] \prod_{l=1}^k \left(\frac{\lambda_i}{\lambda_{j_l}} - 1 \right) \right| \\
&\quad + \frac{1}{(D_k \hat{D}_k)} \left| \sum_{(j_1, \dots, j_k) \in I_k^+} \hat{D}_{j_1, \dots, j_k} \prod_{l=1}^k \left(\frac{\lambda_i}{\lambda_{j_l}} - 1 \right) \right| \left| \sum_{(j_1, \dots, j_k) \in I_k^+} \left[D_{j_1, \dots, j_k} - \hat{D}_{j_1, \dots, j_k} \right] \right|
\end{aligned}$$

Besides we have

$$\begin{aligned}
\left| \sum_{(j_1, \dots, j_k) \in I_k^+} \hat{D}_{j_1, \dots, j_k} \prod_{l=1}^k \left(\frac{\lambda_i}{\lambda_{j_l}} - 1 \right) \right| &\leq \sum_{(j_1, \dots, j_k) \in I_k^+} \hat{D}_{j_1, \dots, j_k} \left[\max_{I_k^+} \left(\prod_{l=1}^k \left| \frac{\lambda_i}{\lambda_{j_l}} - 1 \right| \right) \right], \\
\left| \sum_{(j_1, \dots, j_k) \in I_k^+} \left[D_{j_1, \dots, j_k} - \hat{D}_{j_1, \dots, j_k} \right] \prod_{l=1}^k \left(\frac{\lambda_i}{\lambda_{j_l}} - 1 \right) \right| &= \left| \sum_{(j_1, \dots, j_k) \in I_k^+} D_{j_1, \dots, j_k} \left(1 - \frac{\hat{p}_{j_1}^2 \dots \hat{p}_{j_k}^2}{p_{j_1}^2 \dots p_{j_k}^2} \right) \prod_{l=1}^k \left(\frac{\lambda_i}{\lambda_{j_l}} - 1 \right) \right| \\
&\leq \sum_{(j_1, \dots, j_k) \in I_k^+} D_{j_1, \dots, j_k} \left[\max_{I_k^+} \left(1 - \frac{\hat{p}_{j_1}^2 \dots \hat{p}_{j_k}^2}{p_{j_1}^2 \dots p_{j_k}^2} \right) \right] \left[\max_{I_k^+} \left(\prod_{l=1}^k \left| \frac{\lambda_i}{\lambda_{j_l}} - 1 \right| \right) \right]
\end{aligned}$$

and

$$\begin{aligned}
\left| \sum_{(j_1, \dots, j_k) \in I_k^+} \left[D_{j_1, \dots, j_k} - \hat{D}_{j_1, \dots, j_k} \right] \right| &= \left| \sum_{(j_1, \dots, j_k) \in I_k^+} D_{j_1, \dots, j_k} \left(1 - \frac{\hat{p}_{j_1}^2 \dots \hat{p}_{j_k}^2}{p_{j_1}^2 \dots p_{j_k}^2} \right) \right| \\
&\leq \sum_{(j_1, \dots, j_k) \in I_k^+} D_{j_1, \dots, j_k} \left[\max_{I_k^+} \left(1 - \frac{\hat{p}_{j_1}^2 \dots \hat{p}_{j_k}^2}{p_{j_1}^2 \dots p_{j_k}^2} \right) \right].
\end{aligned}$$

Therefore we get

$$\left| \hat{Q}_k(\lambda_i) - Q_k^*(\lambda_i) \right| \leq 2 \left[\max_{I_k^+} \left(1 - \frac{\hat{p}_{j_1}^2 \dots \hat{p}_{j_k}^2}{p_{j_1}^2 \dots p_{j_k}^2} \right) \right] \left[\max_{I_k^+} \left(\prod_{l=1}^k \left| \frac{\lambda_i}{\lambda_{j_l}} - 1 \right| \right) \right] \quad (23)$$

where

$$\frac{\hat{p}_{j_1}^2 \dots \hat{p}_{j_k}^2}{p_{j_1}^2 \dots p_{j_k}^2} = \left(1 + \frac{\varepsilon_{j_1}}{p_{j_1}} \right)^2 \dots \left(1 + \frac{\varepsilon_{j_k}}{p_{j_k}} \right)^2.$$

From Proposition 7.3 and (H.2) we have that there exists a constant $C > 1$ such that

$$\left(1 - \frac{C}{\sqrt{L}} \sqrt{\frac{\log n}{n}}\right)^{2k} \leq \frac{\hat{p}_{j_1}^2 \dots \hat{p}_{j_k}^2}{p_{j_1}^2 \dots p_{j_k}^2} \leq \left(1 + \frac{C}{\sqrt{L}} \sqrt{\frac{\log n}{n}}\right)^{2k} \quad (24)$$

with probability at least $1 - n^{1-C}$.

From (23) and (24) we deduce that there exists a constant \tilde{C} such that with probability at least $1 - n^{1-C}$ where $C > 1$,

$$\left| \hat{Q}_k(\lambda_i) - Q_k^*(\lambda_i) \right| \leq \frac{2k\tilde{C}}{\sqrt{L}} \sqrt{\frac{\log n}{n}} \left[\max_{I_k^+} \left(\prod_{l=1}^k \left| \frac{\lambda_i}{\lambda_{j_l}} - 1 \right| \right) \right]. \quad (25)$$

Finally using again Proposition 7.3 and (25) we get

$$\begin{aligned} \frac{1}{n} \left\| \left(\hat{Q}_k(X^T X) - Q_k^*(X^T X) \right) Y \right\|^2 &= \frac{1}{n} \sum_{i=1}^n \left(\hat{Q}_k(\lambda_i) - Q_k^*(\lambda_i) \right)^2 \frac{\hat{p}_i^2}{p_i^2} \\ &\leq \frac{4k^2\tilde{C}^2 \log n}{L n^2} \left(1 + \frac{C}{\sqrt{L}} \sqrt{\frac{\log n}{n}} \right) \sum_{i=1}^n \left[\max_{I_k^+} \left(\prod_{l=1}^k \left| \frac{\lambda_i}{\lambda_{j_l}} - 1 \right| \right)^2 p_i^2 \right] \end{aligned}$$

where

$$\sum_{i=1}^n \left[\max_{I_k^+} \left(\prod_{l=1}^k \left| \frac{\lambda_i}{\lambda_{j_l}} - 1 \right| \right)^2 p_i^2 \right] = \|X\beta^*\|_W^2$$

with $W = \text{diag}_{1 \leq i \leq n} \left(\max_{I_k^+} \left(\prod_{l=1}^k \left| \frac{\lambda_i}{\lambda_{j_l}} - 1 \right|^2 \right) \right)$.

Remark We can state sharper bounds for the ratio $\frac{\hat{Q}_k(\lambda_i)}{Q_k^*(\lambda_i)}$ for $i = 1$ and $i = n$. Indeed we have

$$\begin{aligned} \frac{\hat{Q}_k(\lambda_i)}{Q_k^*(\lambda_i)} &= \frac{\sum_{(j_1, \dots, j_k) \in I_k^+} \hat{D}_{j_1, \dots, j_k} \prod_{l=1}^k \left(\frac{\lambda_i}{\lambda_{j_l}} - 1 \right) \sum_{(j_1, \dots, j_k) \in I_k^+} D_{j_1, \dots, j_k}}{\sum_{(j_1, \dots, j_k) \in I_k^+} D_{j_1, \dots, j_k} \prod_{l=1}^k \left(\frac{\lambda_i}{\lambda_{j_l}} - 1 \right) \sum_{(j_1, \dots, j_k) \in I_k^+} \hat{D}_{j_1, \dots, j_k}} \\ &= \frac{\sum_{(j_1, \dots, j_k) \in I_k^+} D_{j_1, \dots, j_k} \prod_{l=1}^k \left(\frac{\lambda_i}{\lambda_{j_l}} - 1 \right) \left(\frac{\hat{p}_{j_1}^2 \dots \hat{p}_{j_k}^2}{p_{j_1}^2 \dots p_{j_k}^2} \right)}{\sum_{(j_1, \dots, j_k) \in I_k^+} D_{j_1, \dots, j_k} \prod_{l=1}^k \left(\frac{\lambda_i}{\lambda_{j_l}} - 1 \right)} \frac{\sum_{(j_1, \dots, j_k) \in I_k^+} D_{j_1, \dots, j_k}}{\sum_{(j_1, \dots, j_k) \in I_k^+} D_{j_1, \dots, j_k} \left(\frac{\hat{p}_{j_1}^2 \dots \hat{p}_{j_k}^2}{p_{j_1}^2 \dots p_{j_k}^2} \right)} \quad (26) \end{aligned}$$

We are going to use again concentration inequalities to bound by above the two factors of the product in (26). In fact on the event \mathcal{A} we have (see (24))

$$\left(1 - C \sqrt{\frac{\log n}{nL}} \right)^{2k} \leq \frac{\hat{p}_{j_1}^2 \dots \hat{p}_{j_k}^2}{p_{j_1}^2 \dots p_{j_k}^2}.$$

Therefore, because all the terms D_{j_1, \dots, j_k} are positive, we deduce that

$$\frac{\sum_{(j_1, \dots, j_k) \in I_k^+} D_{j_1, \dots, j_k}}{\sum_{(j_1, \dots, j_k) \in I_k} D_{j_1, \dots, j_k} \left(\frac{\hat{p}_{j_1}^2 \dots \hat{p}_{j_k}^2}{p_{j_1}^2 \dots p_{j_k}^2} \right)} \leq \frac{1}{\left(1 - C \sqrt{\frac{\log n}{nL}} \right)^{2k}}. \quad (27)$$

If we assume $C\sqrt{\frac{\log n}{nL}} < 1$, we get

$$\frac{\sum_{(j_1, \dots, j_k) \in I_k^+} D_{j_1, \dots, j_k}}{\sum_{(j_1, \dots, j_k) \in I_k} D_{j_1, \dots, j_k} \left(\frac{\hat{p}_{j_1}^2 \dots \hat{p}_{j_k}^2}{p_{j_1}^2 \dots p_{j_k}^2} \right)} \leq 1 + O\left(2k\sqrt{\frac{\log n}{nL}}\right) \quad (28)$$

Now we are going to bound by above the first factor in (26) for $i = 1$ and $i = n$. Let $i = 1$. Then for all $(j_1, \dots, j_k) \in I_k^+$ we have $\prod_{l=1}^k (\frac{\lambda_1}{\lambda_{j_l}} - 1) > 0$ and thus all the terms in the first factor are positive. Therefore on the event \mathcal{A} we get

$$\left| \frac{\sum_{(j_1, \dots, j_k) \in I_k^+} D_{j_1, \dots, j_k} \prod_{l=1}^k (\frac{\lambda_1}{\lambda_{j_l}} - 1) \left(\frac{\hat{p}_{j_1}^2 \dots \hat{p}_{j_k}^2}{p_{j_1}^2 \dots p_{j_k}^2} \right)}{\sum_{(j_1, \dots, j_k) \in I_k^+} D_{j_1, \dots, j_k} \prod_{l=1}^k (\frac{\lambda_1}{\lambda_{j_l}} - 1)} \right| \leq \left(1 + C\sqrt{\frac{\log n}{nL}} \right)^{2k} = 1 + O\left(2k\sqrt{\frac{\log n}{nL}}\right). \quad (29)$$

Let $i = n$. Then for all $(j_1, \dots, j_k) \in I_{k,n}^+$ we have $\prod_{l=1}^k (\frac{\lambda_1}{\lambda_{j_l}} - 1) > 0$ if k is even and $\prod_{l=1}^k (\frac{\lambda_1}{\lambda_{j_l}} - 1) < 0$ if k is odd. Thus on the event \mathcal{A} we have

$$\left| \frac{\sum_{(j_1, \dots, j_k) \in I_k^+} \left(D_{j_1, \dots, j_k} \prod_{l=1}^k (\frac{\lambda_n}{\lambda_{j_l}} - 1) \frac{\hat{p}_{j_1}^2 \dots \hat{p}_{j_k}^2}{p_{j_1}^2 \dots p_{j_k}^2} \right)}{\sum_{(j_1, \dots, j_k) \in I_k^+} D_{j_1, \dots, j_k} \prod_{l=1}^k (\frac{\lambda_1}{\lambda_{j_l}} - 1)} \right| \leq \left(1 + C\sqrt{\frac{\log n}{nL}} \right)^{2k} = 1 + O\left(2k\sqrt{\frac{\log n}{nL}}\right). \quad (30)$$

To conclude from (29), (30) and (28) we have for $i = 1$ and $i = n$

$$\left| \frac{\hat{Q}_k(\lambda_i)}{Q_k^*(\lambda_i)} \right| \leq 1 + O\left(2k\sqrt{\frac{\log n}{nL}}\right).$$

□

References

- Blanchard, G. and Mathé, P. (2012). Discrepancy principle for statistical inverse problems with application to conjugate gradient iteration. *Inverse Problems*, 28(11):115011.
- Boulesteix, A.-L. and Strimmer, K. (2007). Partial least squares: a versatile tool for the analysis of high-dimensional genomic data. *Briefings in bioinformatics*, 8(1):32–44.
- Butler, N. A. and Denham, M. C. (2000). The peculiar shrinkage properties of partial least squares regression. *Journal of the Royal Statistical Society: Series B (Statistical Methodology)*, 62(3):585–593.
- De Jong, S. (1995). Pls shrinks. *Journal of chemometrics*, 9(4):323–326.
- Delaigle, A. and Hall, P. (2012). Methodology and theory for partial least squares applied to functional data. *The Annals of Statistics*, 40(1):322–352.
- Engl, H. W., Hanke, M., and Neubauer, A. (1996). *Regularization of inverse problems*, volume 375 of *Mathematics and its Applications*. Kluwer Academic Publishers Group, Dordrecht.

- Frank, I. E. and Friedman, J. H. (1993). A statistical view of some chemometrics regression tools. *Technometrics*, 35(2):109–135.
- Garthwaite, P. H. (1994). An interpretation of partial least squares. *Journal of the American Statistical Association*, 89(425):122–127.
- Goutis, C. (1996). Partial least squares algorithm yields shrinkage estimators. *The Annals of Statistics*, 24(2):816–824.
- Helland, I. S. (1988). On the structure of partial least squares regression. *Communications in statistics-Simulation and Computation*, 17(2):581–607.
- Helland, I. S. (1990). Partial least squares regression and statistical models. *Scandinavian Journal of Statistics*, pages 97–114.
- Helland, I. S. (2001). Some theoretical aspects of partial least squares regression. *Chemometrics and Intelligent Laboratory Systems*, 58(2):97–107.
- Jolliffe, I. T. (1982). A note on the use of principal components in regression. *Applied Statistics*, pages 300–303.
- Krämer, N. (2007). An overview on the shrinkage properties of partial least squares regression. *Computational Statistics*, 22(2):249–273.
- Lê Cao, K.-A., Rossouw, D., Robert-Granié, C., and Besse, P. (2008). A sparse PLS for variable selection when integrating omics data. *Stat. Appl. Genet. Mol. Biol.*, 7(1):Art. 35, 31.
- Lingjaerde, O. C. and Christophersen, N. (2000). Shrinkage structure of partial least squares. *Scandinavian Journal of Statistics*, 27(3):459–473.
- Martens, H. and Naes, T. (1992). *Multivariate calibration*. Wiley.
- Naes, T. and Martens, H. (1985). Comparison of prediction methods for multicollinear data. *Communications in Statistics-Simulation and Computation*, 14(3):545–576.
- Phatak, A. and de Hoog, F. (2002). Exploiting the connection between pls, lanczos methods and conjugate gradients: alternative proofs of some properties of pls. *Journal of Chemometrics*, 16(7):361–367.
- Saad, Y. (1992). *Numerical methods for large eigenvalue problems*, volume 158. SIAM.
- Wold, H. (1985). Partial least squares. *Encyclopedia of statistical sciences*.
- Wold, S., Martens, H., and Wold, H. (1983). The multivariate calibration problem in chemistry solved by the pls method. In *Matrix pencils*, pages 286–293. Springer.
- Wold, S., Sjöström, M., and Eriksson, L. (2001). Pls-regression: a basic tool of chemometrics. *Chemometrics and intelligent laboratory systems*, 58(2):109–130.